

Reproducing Kernel Banach Spaces with the ℓ^1 Norm*

Guohui Song[†] Haizhang Zhang[‡] and Fred J. Hickernell[§]

Abstract

Targeting at sparse learning, we construct Banach spaces \mathcal{B} of functions on an input space X with the following properties: (1) \mathcal{B} possesses an ℓ^1 norm in the sense that \mathcal{B} is isometrically isomorphic to the Banach space of integrable functions on X with respect to the counting measure; (2) point evaluations are continuous linear functionals on \mathcal{B} and are representable through a bilinear form with a kernel function; and (3) regularized learning schemes on \mathcal{B} satisfy the linear representer theorem. Examples of kernel functions admissible for the construction of such spaces are given.

Keywords: reproducing kernel Banach spaces, sparse learning, lasso, basis pursuit, regularization, the representer theorem, the Brownian bridge kernel, the exponential kernel.

1 Introduction

It is now widely known that minimizing a loss function regularized by the ℓ^1 norm yields sparsity in the resulting minimizer. The sparsity is essential for extracting relatively low dimensional features from sample data that usually live in a high dimensional space. When the square loss function is used in regression, the method is known as the lasso in statistics [26]. Recently, the methodology has been applied to compressive sensing where it is referred to as basis pursuit [4, 5]. The purpose of this paper is to establish an appropriate foundation for developing ℓ^1 regularization for machine learning with reproducing kernels.

Past research on learning with kernels [6, 7, 9, 22, 23, 24, 27] has mainly been built upon the theory of reproducing kernel Hilbert spaces (RKHS) [2]. There are many reasons that account for the success from such a choice. RKHS are by definition the Hilbert space of functions where point evaluations are continuous linear functionals. Sample data available for learning are usually modeled by point evaluations of the unknown target function. Therefore, RKHS is a class of function spaces where sampling is stable, a desirable feature in applications. By the Riesz representation theorem, continuous linear functionals on a Hilbert space are representable by the inner product on the space. This gives rise to the representation of point evaluation functionals

*Supported by Guangdong Provincial Government of China through the “Computational Science Innovative Research Team” program.

[†]School of Mathematical and Statistical Sciences, Arizona State University, Tempe, AZ 85287, USA. E-mail address: *gsong9@asu.edu*.

[‡]School of Mathematics and Computational Science and Guangdong Province Key Laboratory of Computational Science, Sun Yat-sen University, Guangzhou 510275, P. R. China. E-mail address: *zhhaizh2@sysu.edu.cn*.

[§]Department of Applied Mathematics, Illinois Institute of Technology, 10 W. 32nd St., Chicago, IL, 60616, USA. E-mail address: *hickernell@iit.edu*. This author’s work was supported in part by National Science Foundation grants DMS-0713848 and DMS-1115392.

on an RKHS by its associated reproducing kernel and leads to the celebrated representer theorem [14] in machine learning. This theorem states that the original minimization problem in a typically infinite dimensional RKHS can be converted into a problem of determining finitely many coefficients in a linear combination of the kernel function with one argument evaluated at the data sites.

For this representer theorem, the nonzero coefficients to be found are generally as many as the sampling points. For the sake of economy, it is hence desirable to regularize the class of candidate functions by some ℓ^1 norm to force most of the coefficients to be zero. An attempt in this direction is the linear programming approach to coefficient based regularization for machine learning [22]. The method lacks a general mathematical foundation like the RKHS though. In particular, it is unknown whether the algorithm results by some representer theorem from a minimization on an infinite dimensional Banach space. A consequence is that the hypothesis error in the learning rate estimate will not go away automatically as in the RKHS case [30].

We aim at combining the reproducing kernel methods and the ℓ^1 regularization technique. Specifically, we desire to construct function spaces with the following properties:

- point evaluation functionals on the space are continuous and can be represented by some kernel function;
- the space possesses an ℓ^1 norm;
- a linear representer theorem holds for regularized learning schemes on the space.

There are three ways of representing continuous point evaluation functionals in a function space: by an inner product, by a semi-inner product [11, 15], or by a bilinear form on the tensor product of the space and its dual space. Since the space we constructed is expected to have an ℓ^1 norm, it can not have an inner product. Semi-inner products are a natural substitute for inner products in Banach spaces. A notion of reproducing kernel Banach spaces (RKBS) was established in [31, 32] via the semi-inner product. The spaces considered there are uniformly convex and uniformly Fréchet differentiable to ensure that continuous linear functionals have a unique representation by the semi-inner product. An infinite dimensional Banach space with the ℓ^1 norm is non-reflexive. As a consequence, there is no guarantee [13] that the semi-inner product is able to represent all continuous point evaluation functionals in such a space. For these reasons, we shall pursue the third approach in this study, that is, to represent the point evaluation functionals by a bilinear form. We briefly introduce the construction and main results of the paper below.

Let X be a prescribed set that we call the input space. The construction starts directly with a complex-valued function K on $X \times X$, which is not necessarily Hermitian. For the constructed space to have the three desirable properties described above, K needs to be an admissible kernel. To introduce this class of functions crucial to our construction, we denote for any set Ω by $\ell^1(\Omega)$ the Banach space of functions on Ω that is integrable with respect to the counting measure on Ω . In other words,

$$\ell^1(\Omega) := \{\mathbf{c} = (c_t \in \mathbb{C} : t \in \Omega) : \|\mathbf{c}\|_{\ell^1(\Omega)} := \sum_{t \in \Omega} |c_t| < +\infty\}.$$

Note that Ω might be uncountable but for every $\mathbf{c} \in \ell^1(\Omega)$, $\text{supp } \mathbf{c} := \{t \in \Omega : c_t \neq 0\}$ must be countable. Finally, we define the set $\mathbb{N}_n := \{1, 2, \dots, n\}$ for all $n \in \mathbb{N}$.

Definition 1.1. *A function K on $X \times X$ is called an admissible kernel for the construction of RKBS on X with the ℓ^1 norm if the following requirements are satisfied:*

(A1) for all sequences $\mathbf{x} = \{x_j : j \in \mathbb{N}_n\} \subseteq X$ of pairwise distinct sampling points, the matrix

$$K[\mathbf{x}] := [K(x_k, x_j) : j, k \in \mathbb{N}_n] \in \mathbb{C}^{n \times n} \quad (1.1)$$

is nonsingular,

(A2) K is bounded, namely, $|K(s, t)| \leq M$ for some positive constant M and all $s, t \in X$,

(A3) for all pairwise distinct $x_j \in X$, $j \in \mathbb{N}$ and $\mathbf{c} \in \ell^1(\mathbb{N})$, $\sum_{j=1}^{\infty} c_j K(x_j, x) = 0$ for all $x \in X$ implies $\mathbf{c} = 0$, and

(A4) for all pairwise distinct $x_1, x_2, \dots, x_{n+1} \in X$,

$$\|(K[\mathbf{x}])^{-1} K_{\mathbf{x}}(x_{n+1})\|_{\ell^1(\mathbb{N}_n)} \leq 1, \quad (1.2)$$

where $K_{\mathbf{x}}(x) = (K(x, x_j) : j \in \mathbb{N}_n)^T \in \mathbb{C}^n$.

The following theorem will be proved in the next three sections.

Theorem 1.2. *If K is an admissible kernel on $X \times X$ then*

$$\mathcal{B} := \left\{ \sum_{t \in \text{supp } \mathbf{c}} c_t K(t, \cdot) : \mathbf{c} \in \ell^1(X) \right\} \text{ with the norm } \left\| \sum_{t \in \text{supp } \mathbf{c}} c_t K(t, \cdot) \right\|_{\mathcal{B}} := \|\mathbf{c}\|_{\ell^1(X)} \quad (1.3)$$

and \mathcal{B}^\sharp , the completion of the vector space of functions $\sum_{j=1}^n c_j K(\cdot, x_j)$, $x_j \in X$ under the supremum norm

$$\left\| \sum_{j=1}^n c_j K(\cdot, x_j) \right\|_{\mathcal{B}^\sharp} := \sup \left\{ \left| \sum_{j=1}^n c_j K(x, x_j) \right| : x \in X \right\},$$

are both Banach spaces of functions on X where point evaluations are continuous linear functionals. In addition, the bilinear form

$$\left\langle \sum_{j=1}^n a_j K(s_j, \cdot), \sum_{k=1}^m b_k K(\cdot, t_k) \right\rangle_K := \sum_{j=1}^n \sum_{k=1}^m a_j b_k K(s_j, t_k), \quad s_j, t_k \in X \quad (1.4)$$

can be extended to $\mathcal{B} \times \mathcal{B}^\sharp$ such that

$$|\langle f, g \rangle_K| \leq \|f\|_{\mathcal{B}} \|g\|_{\mathcal{B}^\sharp} \text{ for all } f \in \mathcal{B}, g \in \mathcal{B}^\sharp$$

and

$$\langle f, K(\cdot, x) \rangle_K = f(x), \quad \langle K(x, \cdot), g \rangle_K = g(x) \text{ for all } x \in X, f \in \mathcal{B}, g \in \mathcal{B}^\sharp.$$

Furthermore, for every regularized learning scheme of the form

$$\inf_{f \in \mathcal{B}} V(f(x_1), f(x_2), \dots, f(x_n)) + \mu \phi(\|f\|_{\mathcal{B}}),$$

where μ is a positive regularization parameter, V and ϕ are nonnegative continuous functions with $\lim_{t \rightarrow \infty} \phi(t) = +\infty$, there exists a minimizer, f_0 , of the form

$$f_0(x) = \sum_{j=1}^n c_j K(x_j, x), \quad x \in X$$

for some coefficients $c_j \in \mathbb{C}$, $j \in \mathbb{N}_n$.

Conversely, for the constructed spaces \mathcal{B} and \mathcal{B}^\sharp to enjoy those desirable properties, K must be an admissible kernel on $X \times X$.

The organization of the paper is as follows. We first present a general construction of Banach spaces of functions with a reproducing kernel in the next section. In Section 3, we specify the construction to the building of RKBS with the ℓ^1 norm as described in Theorem 1.2. In Section 4, we study the conditions on the reproducing kernel so that regularized learning schemes on the constructed spaces satisfy the linear representer theorem. In the last section, we show that the Brownian bridge kernel and the exponential kernel are admissible kernels. In the final section, condition (A4), the most stringent condition in Definition 1.1 is relaxed, which leads to a modified version of Theorem 1.2.

2 A General Construction

To ensure that there exists a reproducing kernel, we shall start the construction of the Banach space based on such a function. Let X be an input space and let K be a function on $X \times X$. Introduce the vector space

$$\mathcal{B}_0 := \text{span} \{K(x, \cdot) : x \in X\}.$$

Note that unlike reproducing kernels for Hilbert spaces, this K is not necessarily symmetric in its arguments or positive definite. Suppose that a norm $\|\cdot\|_{\mathcal{B}_0}$ is imposed on \mathcal{B}_0 such that point evaluation functionals are continuous on \mathcal{B}_0 . That is, for any $x \in X$, there exists a positive constant M_x such that

$$|\delta_x(f)| = |f(x)| \leq M_x \|f\|_{\mathcal{B}_0} \text{ for all } f \in \mathcal{B}_0. \quad (2.1)$$

The function K and the norm on \mathcal{B}_0 will be explicitly given in a specific construction.

In [31, 33, 32], a vector space \mathcal{B} is called an RKBS on X if it is a uniformly convex and uniformly Fréchet differentiable Banach space of functions on X and point evaluation functionals are continuous on \mathcal{B} . The uniform convexity and uniform Fréchet differentiability were imposed there to ensure the existence of a reproducing kernel for representing the point evaluation functionals. By the results to be established in the current paper, these stronger conditions are not necessary. To accommodate the search for alternatives, we introduce the following definitions.

Definition 2.1. *The space \mathcal{B} is called a Banach space of functions if the point evaluation functionals are consistent with the norm on \mathcal{B} in the sense that for all $f \in \mathcal{B}$, $\|f\|_{\mathcal{B}} = 0$ if and only if f vanishes everywhere on X . A Banach space \mathcal{B} of functions on X is said to be a pre-RKBS on X if point evaluations are continuous linear functionals on \mathcal{B} .*

We plan to complete \mathcal{B}_0 by the norm $\|\cdot\|_{\mathcal{B}_0}$ to obtain a pre-RKBS \mathcal{B} . Two things need to be checked for the approach to succeed. An abstract completion of \mathcal{B}_0 might not consist of functions, or might not have bounded point evaluation functionals. We shall present a Banach completion process that yields a space of functions. Let $\{f_n : n \in \mathbb{N}\}$ be a Cauchy sequence in \mathcal{B}_0 . Since point evaluation functionals are continuous on \mathcal{B}_0 , for any $x \in X$, the sequence $\{f_n(x) : n \in \mathbb{N}\}$ converges in \mathbb{C} . We denote the limit by $f(x)$, which defines a function on X . One sees that two equivalent Cauchy sequences in \mathcal{B}_0 give the same function. We let \mathcal{B} be composed of all such limit functions with the norm $\|f\|_{\mathcal{B}} := \lim_{n \rightarrow \infty} \|f_n\|_{\mathcal{B}_0}$.

To investigate conditions for \mathcal{B} to be a pre-RKBS, we need to invoke the following assumption.

Definition 2.2. *A normed vector space V of functions on X satisfies the Norm Consistency Property if for every Cauchy sequence $\{f_n : n \in \mathbb{N}\}$ in V , $\lim_{n \rightarrow \infty} f_n(x) = 0$ for all $x \in X$ implies $\lim_{n \rightarrow \infty} \|f_n\|_V = 0$.*

Proposition 2.3. *The norm $\|\cdot\|_{\mathcal{B}}$ is well-defined and makes \mathcal{B} a pre-RKBS on X if and only if \mathcal{B}_0 satisfies the Norm Consistency Property.*

Proof. We first show the necessity. If \mathcal{B} is a Banach space then $\|\cdot\|_{\mathcal{B}}$ is a well-defined norm. The validity of the Norm Consistency Property follows directly from $\|0\|_{\mathcal{B}} = 0$.

We next prove the sufficiency. Suppose that the Norm Consistency Property holds for \mathcal{B}_0 . We first show that $\|\cdot\|_{\mathcal{B}}$ is a well-defined norm. Suppose that $\{f_n : n \in \mathbb{N}\}$ and $\{g_n : n \in \mathbb{N}\}$ are both Cauchy sequences in \mathcal{B}_0 such that $\lim_{n \rightarrow \infty} f_n(x) = \lim_{n \rightarrow \infty} g_n(x)$ for all $x \in X$. We need to show that $\lim_{n \rightarrow \infty} \|f_n\|_{\mathcal{B}_0} = \lim_{n \rightarrow \infty} \|g_n\|_{\mathcal{B}_0}$. Clearly, $f_n - g_n$ forms a Cauchy sequence in \mathcal{B}_0 . Since $\lim_{n \rightarrow \infty} (f_n - g_n)(x) = 0$ for all $x \in X$, it follows from the Norm Consistency Property that $\lim_{n \rightarrow \infty} \|f_n - g_n\|_{\mathcal{B}_0} = 0$, which implies $\lim_{n \rightarrow \infty} \|f_n\|_{\mathcal{B}_0} = \lim_{n \rightarrow \infty} \|g_n\|_{\mathcal{B}_0}$. Therefore, $\|\cdot\|_{\mathcal{B}}$ is well-defined. As a result, \mathcal{B} is isometrically isomorphic to the abstract Banach space that is the completion of \mathcal{B}_0 . It implies that \mathcal{B} is a Banach space and \mathcal{B}_0 is dense in \mathcal{B} . Moreover, it follows immediately from the Norm Consistency Property that \mathcal{B} is a Banach space of functions. It remains to show that the point evaluation functional δ_x is continuous on \mathcal{B} for all $x \in X$. Let $x \in X$ and $f \in \mathcal{B}$. By definition, there exists a Cauchy sequence $\{f_n : n \in \mathbb{N}\}$ in \mathcal{B}_0 such that

$$f(x) = \lim_{n \rightarrow \infty} f_n(x) \text{ for all } x \in X, \quad \text{and} \quad \|f\|_{\mathcal{B}} = \lim_{n \rightarrow \infty} \|f_n\|_{\mathcal{B}_0}.$$

Since δ_x is continuous on \mathcal{B}_0 , there exists a positive constant M_x such that

$$|f_n(x)| \leq M_x \|f_n\|_{\mathcal{B}_0} \quad \text{for all } n \in \mathbb{N}.$$

Taking the limits on both sides, we have $|f(x)| \leq M_x \|f\|_{\mathcal{B}}$. The proof is complete. \square

In the rest of this section, we assume the Norm Consistency Property for \mathcal{B}_0 and aim at deriving a reproducing kernel for \mathcal{B} . To this end, we set

$$\mathcal{B}_0^\sharp := \text{span} \{K(\cdot, x) : x \in X\}$$

and define a bilinear form $\langle \cdot, \cdot \rangle_K$ on $\mathcal{B}_0 \times \mathcal{B}_0^\sharp$ by (1.4). It is straightforward to observe that

$$\langle f, K(\cdot, x) \rangle_K = f(x), \quad \langle K(x, \cdot), g \rangle_K = g(x) \text{ for all } f \in \mathcal{B}_0, g \in \mathcal{B}_0^\sharp \text{ and } x \in X.$$

It means (1.4) is well defined and that K is able to reproduce the point evaluations of functions on \mathcal{B}_0 via this bilinear form. We need to extend this property to the whole space \mathcal{B} in order to claim that it is a reproducing kernel for \mathcal{B} . For this purpose, we define another norm

$$\|g\|_{\mathcal{B}_0^\sharp} := \sup_{f \in \mathcal{B}_0, f \neq 0} \frac{|\langle f, g \rangle_K|}{\|f\|_{\mathcal{B}_0}}, \quad g \in \mathcal{B}_0^\sharp. \quad (2.2)$$

The next result indicates that the above norm is well-defined.

Proposition 2.4. *The norm $\|\cdot\|_{\mathcal{B}_0^\sharp}$ is well-defined and point evaluation functionals are continuous on \mathcal{B}_0^\sharp if and only if point evaluation functionals are continuous on \mathcal{B}_0 .*

Proof. We begin with the sufficiency. Suppose that point evaluation functionals are continuous on \mathcal{B}_0 . That is, for any $x \in X$ there exists a positive constant M_x satisfying (2.1). Let $g \in \mathcal{B}_0^\sharp$.

It must be of the form $g = \sum_{j=1}^n a_j K(\cdot, x_j)$ for some $a_j \in \mathbb{C}$ and $x_j \in X$, $j \in \mathbb{N}_n$, $n \in \mathbb{N}$. We have for all $f \in \mathcal{B}_0$

$$\frac{|\langle f, g \rangle_K|}{\|f\|_{\mathcal{B}_0}} = \frac{|\langle f, \sum_{j=1}^n a_j K(\cdot, x_j) \rangle_K|}{\|f\|_{\mathcal{B}_0}} = \frac{|\sum_{j=1}^n a_j f(x_j)|}{\|f\|_{\mathcal{B}_0}} \leq \sum_{j=1}^n |a_j| M_{x_j},$$

which implies that $\|g\|_{\mathcal{B}_0^\sharp}$ is well-defined. We next prove that point evaluation functionals are continuous on \mathcal{B}_0^\sharp . By (2.2), we have for all $f \in \mathcal{B}_0, g \in \mathcal{B}_0^\sharp$

$$|\langle f, g \rangle_K| \leq \|f\|_{\mathcal{B}_0} \|g\|_{\mathcal{B}_0^\sharp}. \quad (2.3)$$

For any $x \in X$, taking $f = K(x, \cdot)$ in the above inequality yields that

$$|g(x)| = |\langle K(x, \cdot), g \rangle_K| \leq \|K(x, \cdot)\|_{\mathcal{B}_0} \|g\|_{\mathcal{B}_0^\sharp} \quad \text{for all } g \in \mathcal{B}_0^\sharp.$$

It follows that the point evaluation functional δ_x is continuous on \mathcal{B}_0^\sharp as $\|K(x, \cdot)\|_{\mathcal{B}_0}$ is a constant independent of g .

We next turn to the necessity. Suppose $\|g\|_{\mathcal{B}_0^\sharp}$ is well-defined for all $g \in \mathcal{B}_0^\sharp$. For any $x \in X$, letting $g = K(\cdot, x)$ in (2.3) yields

$$|f(x)| \leq \|K(\cdot, x)\|_{\mathcal{B}_0^\sharp} \|f\|_{\mathcal{B}_0},$$

which implies that point evaluation functionals are continuous on \mathcal{B}_0 . \square

We complete \mathcal{B}_0^\sharp using the norm $\|\cdot\|_{\mathcal{B}_0^\sharp}$ to a Banach space \mathcal{B}^\sharp by the process described before Proposition 2.3. We have the following observation similar to that about the space \mathcal{B} .

Proposition 2.5. *The space \mathcal{B}^\sharp is a pre-RKBS on X if and only if the normed vector space \mathcal{B}_0^\sharp satisfies the Norm Consistency Property.*

In the following discussion, suppose that \mathcal{B}_0^\sharp endowed with the norm $\|\cdot\|_{\mathcal{B}_0^\sharp}$ has the Norm Consistency Property. By applying the Hahn-Banach extension theorem twice, we can extend the bilinear form $\langle \cdot, \cdot \rangle_K$ from $\mathcal{B}_0 \times \mathcal{B}_0^\sharp$ to $\mathcal{B} \times \mathcal{B}^\sharp$ in a unique way such that

$$|\langle f, g \rangle_K| \leq \|f\|_{\mathcal{B}} \|g\|_{\mathcal{B}^\sharp}, \quad f \in \mathcal{B}, \quad g \in \mathcal{B}^\sharp. \quad (2.4)$$

The next result tells that the definition of $\|\cdot\|_{\mathcal{B}_0^\sharp}$ in (2.2) can be extended to \mathcal{B}^\sharp .

Proposition 2.6. *Suppose that point evaluation functionals are continuous on \mathcal{B}_0 . If both \mathcal{B}_0 and \mathcal{B}_0^\sharp satisfy the Norm Consistency Property then we have*

$$\|g\|_{\mathcal{B}^\sharp} = \sup_{f \in \mathcal{B}, f \neq 0} \frac{|\langle f, g \rangle_K|}{\|f\|_{\mathcal{B}}}, \quad g \in \mathcal{B}^\sharp. \quad (2.5)$$

Proof. By (2.4), the right hand side above is bounded by the left hand side. We only need to prove the other direction of the inequality. We first show it for functions in \mathcal{B}_0^\sharp . Let $g \in \mathcal{B}_0^\sharp$. It is straightforward to observe that

$$\|g\|_{\mathcal{B}^\sharp} = \sup_{f \in \mathcal{B}_0, f \neq 0} \frac{|\langle f, g \rangle_K|}{\|f\|_{\mathcal{B}}} \leq \sup_{f \in \mathcal{B}, f \neq 0} \frac{|\langle f, g \rangle_K|}{\|f\|_{\mathcal{B}}}. \quad (2.6)$$

Now let g be an arbitrary but fixed function in \mathcal{B}^\sharp . Since \mathcal{B}_0^\sharp is dense in \mathcal{B}^\sharp , there exists $\{g_n : n \in \mathbb{N}\} \subseteq \mathcal{B}_0^\sharp$ such that $\|g - g_n\|_{\mathcal{B}^\sharp} \rightarrow 0$ as $n \rightarrow \infty$. This together with (2.6) implies

$$\|g\|_{\mathcal{B}^\sharp} = \lim_{n \rightarrow \infty} \|g_n\|_{\mathcal{B}^\sharp} \leq \lim_{n \rightarrow \infty} \sup_{f \in \mathcal{B}, f \neq 0} \frac{|\langle f, g_n \rangle_K|}{\|f\|_{\mathcal{B}}}.$$

Note that

$$\frac{|\langle f, g_n \rangle_K|}{\|f\|_{\mathcal{B}}} \leq \frac{|\langle f, g \rangle_K|}{\|f\|_{\mathcal{B}}} + \frac{|\langle f, g - g_n \rangle_K|}{\|f\|_{\mathcal{B}}} \leq \frac{|\langle f, g \rangle_K|}{\|f\|_{\mathcal{B}}} + \|g - g_n\|_{\mathcal{B}^\sharp}.$$

It follows from the above two equations that

$$\|g\|_{\mathcal{B}^\sharp} \leq \lim_{n \rightarrow \infty} \sup_{f \in \mathcal{B}, f \neq 0} \left[\frac{|\langle f, g \rangle_K|}{\|f\|_{\mathcal{B}}} + \|g - g_n\|_{\mathcal{B}^\sharp} \right] = \sup_{f \in \mathcal{B}, f \neq 0} \frac{|\langle f, g \rangle_K|}{\|f\|_{\mathcal{B}}},$$

which completes the proof. \square

We next present necessary and sufficient conditions for K to be able to reproduce point evaluation functionals on \mathcal{B} and \mathcal{B}^\sharp by the bilinear form. We shall see that assuming the Norm Consistency Property, both \mathcal{B} and \mathcal{B}^\sharp are Banach spaces of functions on X such that the point evaluation functionals are continuous and can be represented by the bilinear form with the function K . It is in this sense that \mathcal{B} and \mathcal{B}^\sharp are said to be a *reproducing kernel Banach space with the reproducing kernel K* .

Theorem 2.7. *Suppose that \mathcal{B}_0 and \mathcal{B}_0^\sharp satisfy the Norm Consistency Property. Then both \mathcal{B} and \mathcal{B}^\sharp are pre-RKBS on X and the kernel K reproduces function values via the bilinear form, namely,*

$$\langle f, K(\cdot, x) \rangle_K = f(x) \text{ for all } x \in X \text{ and } f \in \mathcal{B} \quad (2.7)$$

and

$$\langle K(x, \cdot), g \rangle_K = g(x) \text{ for all } x \in X \text{ and } g \in \mathcal{B}^\sharp. \quad (2.8)$$

Thus, \mathcal{B} and \mathcal{B}^\sharp are reproducing kernel Banach spaces (RKBS).

Proof. By Propositions 2.3 and 2.5, both \mathcal{B} and \mathcal{B}^\sharp are pre-RKBS on X . For each $f \in \mathcal{B}$, there exists a sequence $\{f_n : n \in \mathbb{N}\} \subseteq \mathcal{B}_0$ convergent to f . As a consequence, we have for any $x \in X$

$$f(x) = \lim_{n \rightarrow \infty} f_n(x) = \lim_{n \rightarrow \infty} \langle f_n, K(\cdot, x) \rangle_K.$$

By (2.4), $\langle \cdot, K(\cdot, x) \rangle_K$ is a bounded linear functional on \mathcal{B} , which implies

$$\lim_{n \rightarrow \infty} \langle f_n, K(\cdot, x) \rangle_K = \langle f, K(\cdot, x) \rangle_K.$$

Combining the above two equations proves (2.7). Equation (2.8) can be proved similarly. \square

We next discuss the relationship between the space \mathcal{B}^\sharp and the dual space \mathcal{B}^* of \mathcal{B} . It is clear by (2.4) and (2.5) that the mapping \mathcal{L} from \mathcal{B}^\sharp to \mathcal{B}^* defined by the bilinear form,

$$(\mathcal{L}g)(f) := \langle f, g \rangle_K, \quad f \in \mathcal{B}, \quad g \in \mathcal{B}^\sharp, \quad (2.9)$$

is isometric and linear. In other words, \mathcal{L} is an embedding from \mathcal{B}^\sharp to \mathcal{B}^* . We next present a necessary and sufficient condition for it to be surjective.

Proposition 2.8. *Suppose that both \mathcal{B}_0 and \mathcal{B}_0^\sharp satisfy the Norm Consistency Property. The mapping \mathcal{L} defined by (2.9) is surjective onto \mathcal{B}^* if and only if for any proper closed subspace $\mathcal{M} \subsetneq \mathcal{B}$, the orthogonal space $\mathcal{M}^\perp := \{g \in \mathcal{B}^\sharp : \langle f, g \rangle_K = 0 \text{ for all } f \in \mathcal{M}\}$ is nontrivial.*

Proof. We first prove the necessity. For any proper closed subspace $\mathcal{M} \subsetneq \mathcal{B}$, by the Hahn-Banach theorem, there exists a nontrivial functional $\nu \in \mathcal{B}^*$ such that $\nu(f) = 0$ for all $f \in \mathcal{M}$. If \mathcal{L} is surjective then there exists a function $g \in \mathcal{B}^\sharp$ such that $\mathcal{L}(g) = \nu$, namely, $\nu(f) = \langle f, g \rangle_K$ for all $f \in \mathcal{B}$. It follows that $g \in \mathcal{M}^\perp$ and $g \neq 0$ as ν is nontrivial.

We next show the sufficiency. Let ν be a nontrivial functional in \mathcal{B}^* . Then its kernel $\ker(\nu)$ is a proper closed subspace of \mathcal{B} . By assumption, there exists a nonzero function $g \in \mathcal{M}^\perp$. This enables us to find a function $f_0 \in \mathcal{B} \setminus \mathcal{M}$ such that $\langle f_0, g \rangle_K \neq 0$ and $\nu(f_0) = 1$. Set $g_0 := g / \langle f_0, g \rangle_K$. Since $f - \nu(f)f_0 \in \ker(\nu)$ for all $f \in \mathcal{B}$, we get for any $f \in \mathcal{M}$

$$\langle f, g_0 \rangle_K = \langle f - \nu(f)f_0, g_0 \rangle_K + \langle \nu(f)f_0, g_0 \rangle_K = \nu(f)\langle f_0, g_0 \rangle_K = \nu(f),$$

which implies that \mathcal{L} is surjective. \square

We close the section with a conclusion on the general construction and the related results presented above.

Theorem 2.9. *Suppose that*

- (a) *the vector space $\mathcal{B}_0 = \text{span}\{K(x, \cdot) : x \in X\}$ with the norm $\|\cdot\|_{\mathcal{B}_0}$ has the Norm Consistency Property, and*
- (b) *point evaluation functionals are continuous on \mathcal{B}_0 .*

Then the following statements hold true:

- (1) *\mathcal{B}_0 can be completed to a pre-RKBS \mathcal{B} on X ;*
- (2) *the norm $\|\cdot\|_{\mathcal{B}_0^\sharp}$ given by (2.2) is well-defined and point evaluation functionals are bounded on \mathcal{B}_0^\sharp with respect to this norm;*
- (3) *if \mathcal{B}_0^\sharp satisfies the Norm Consistency Property as well then \mathcal{B}_0^\sharp can be completed to an RKBS \mathcal{B}^\sharp and K is the reproducing kernel for both \mathcal{B} and \mathcal{B}^\sharp in the sense that (2.7) and (2.8) hold true. In this case, \mathcal{B}^\sharp can be isometrically embedded into \mathcal{B}^* via the bilinear form, and the embedding is surjective if and only if for any proper closed subspace \mathcal{M} of \mathcal{B} , \mathcal{M}^\perp is nontrivial.*

3 RKBS with the ℓ^1 Norm

We shall follow the procedures in Theorem 2.9 to construct an RKBS with the ℓ^1 norm in this section. To start, we let K be a bounded function on $X \times X$ such that

$$K(x_j, \cdot), j \in \mathbb{N}_n \text{ are linearly independent for all pairwise distinct points } x_j \in X, j \in \mathbb{N}_n. \quad (3.1)$$

Note that this assumption is implied by Admissibility Assumption (A1), but is somewhat weaker than (A1). Introduce an ℓ^1 norm on $\mathcal{B}_0 = \text{span}\{K(x, \cdot) : x \in X\}$ by setting for all finitely many pairwise distinct points $x_j \in X$ and constants $c_j \in \mathbb{C}$, $j \in \mathbb{N}_m$, $m \in \mathbb{N}$

$$\left\| \sum_{j=1}^m c_j K(x_j, \cdot) \right\|_{\mathcal{B}_0} := \sum_{j=1}^m |c_j|. \quad (3.2)$$

Since K is bounded, it is clear that point evaluation functionals are bounded on \mathcal{B}_0 . We next check the important Norm Consistency Property and find that it is implied by the Admissibility Assumption above.

Proposition 3.1. *The space \mathcal{B}_0 with the norm (3.2) satisfies the Norm Consistency Property if and only if K satisfies (A3).*

Proof. We first show the necessity. Suppose that for some $\mathbf{c} \in \ell^1(\mathbb{N})$ and pairwise distinct $\{x_j \in X : j \in \mathbb{N}\}$, $\sum_{j=1}^{\infty} c_j K(x_j, x) = 0$ for all $x \in X$. Let $f_n := \sum_{j=1}^n c_j K(x_j, \cdot)$ for all $n \in \mathbb{N}$. Since $\mathbf{c} \in \ell^1(\mathbb{N})$, $\{f_n : n \in \mathbb{N}\}$ forms a Cauchy sequence in \mathcal{B}_0 . Moreover, $\lim_{n \rightarrow \infty} f_n(x) = 0$ for all $x \in X$ as K is bounded on $X \times X$. It follows from the Norm Consistency Property that $\lim_{n \rightarrow \infty} \|f_n\|_{\mathcal{B}_0} = \lim_{n \rightarrow \infty} \sum_{j=1}^n |c_j| = \|\mathbf{c}\|_{\ell^1(\mathbb{N})} = 0$. Therefore, (A3) holds true.

On the other hand, suppose that K satisfies (A3). Let $\{f_n : n \in \mathbb{N}\}$ be a Cauchy sequence in \mathcal{B}_0 with $\lim_{n \rightarrow \infty} f_n(x) = 0$ for all $x \in X$. We can find pairwise distinct $x_j \in X$, $j \in \mathbb{N}$ such that for any $n \in \mathbb{N}$

$$f_n = \sum_{j=1}^{\infty} c_{n,j} K(x_j, \cdot),$$

where $\mathbf{c}_n := (c_{n,j} : j \in \mathbb{N})$ has finitely many nonzero components. By definition (3.2), $\{\mathbf{c}_n : n \in \mathbb{N}\}$ is a Cauchy sequence in $\ell^1(\mathbb{N})$. Let \mathbf{c} be its limit in $\ell^1(\mathbb{N})$ and define

$$f := \sum_{j=1}^{\infty} c_j K(x_j, \cdot).$$

Suppose that $|K(s, t)| \leq M$ for some positive constant M and all $s, t \in X$. A direct calculation gives that for any $x \in X$

$$|f_n(x) - f(x)| = \left| \sum_{j=1}^{\infty} (c_{n,j} - c_j) K(x_j, x) \right| \leq M \|\mathbf{c}_n - \mathbf{c}\|_{\ell^1(\mathbb{N})}.$$

It follows that $\lim_{n \rightarrow \infty} f_n(x) = f(x)$ for all $x \in X$. Since $\lim_{n \rightarrow \infty} f_n(x) = 0$ for all $x \in X$, we have $f(x) = 0$ for all $x \in X$. By (A3), $\mathbf{c} = 0$, which implies

$$\lim_{n \rightarrow \infty} \|f_n\|_{\mathcal{B}_0} = \lim_{n \rightarrow \infty} \|\mathbf{c}_n\|_{\ell^1(\mathbb{N})} = \|\mathbf{c}\|_{\ell^1(\mathbb{N})} = 0.$$

The proof is complete. \square

Functions K satisfying property (A3) will be given later. We assume for the time being that (A3) holds true. One sees from the proof of Proposition 3.1 that \mathcal{B} has the form (1.3). We remark that in the preparation of the paper, we came across a Banach space with a form similar to (1.3) used in [30] for error estimates with linear programming regularization. One observes from (1.3) that $\ell^1(X)$ is isometrically isomorphic to \mathcal{B} through the mapping

$$\Phi(\mathbf{c}) := \sum_{t \in X} c_t K(t, \cdot), \quad \mathbf{c} \in \ell^1(X).$$

In this sense, we say that \mathcal{B} is a pre-RKBS on X with the ℓ^1 norm. It remains to derive a reproducing kernel for it. By Theorem 2.7, it suffices to check the Norm Consistency Property

for \mathcal{B}_0^\sharp . We shall show that the Norm Consistency Property automatically holds true for \mathcal{B}_0^\sharp without any additional requirement. To this end, we first calculate a specific form of the norm $\|\cdot\|_{\mathcal{B}_0^\sharp}$.

Denote for any function g on X by $\|g\|_{L^\infty(X)}$ the supremum of $|g(x)|$ over $x \in X$.

Lemma 3.2. *There holds for any function $g \in \mathcal{B}_0^\sharp$ that $\|g\|_{\mathcal{B}_0^\sharp} = \|g\|_{L^\infty(X)}$.*

Proof. We first prove that $\|g\|_{\mathcal{B}_0^\sharp}$ is bounded by $\|g\|_{L^\infty(X)}$. Any $f \in \mathcal{B}_0$ has the form $f = \sum_{j=1}^n c_j K(x_j, \cdot)$ for some $c_j \in \mathbb{C}$ and pairwise distinct $x_j \in X$, $j \in \mathbb{N}_n$. We verify that

$$|\langle f, g \rangle_K| = \left| \left\langle \sum_{j=1}^n c_j K(x_j, \cdot), g \right\rangle \right| = \left| \sum_{j=1}^n c_j g(x_j) \right| \leq \|g\|_{L^\infty(X)} \sum_{j=1}^n |c_j| = \|g\|_{L^\infty(X)} \|f\|_{\mathcal{B}_0},$$

which implies $\|g\|_{\mathcal{B}_0^\sharp} \leq \|g\|_{L^\infty(X)}$. For the other direction, we notice for all $x_0 \in X$

$$\|g\|_{\mathcal{B}_0^\sharp} \geq \frac{|\langle K(x_0, \cdot), g \rangle_K|}{\|K(x_0, \cdot)\|_{\mathcal{B}_0}} = |g(x_0)|.$$

Since x_0 is arbitrarily chosen, we have $\|g\|_{\mathcal{B}_0^\sharp} \geq \|g\|_{L^\infty(X)}$. □

We show that the space \mathcal{B}^\sharp is also a pre-RKBS on X .

Lemma 3.3. *The space \mathcal{B}_0^\sharp satisfies the Norm Consistency Property.*

Proof. Let $\{f_n : n \in \mathbb{N}\}$ be a Cauchy sequence in \mathcal{B}_0^\sharp with $\lim_{n \rightarrow \infty} f_n(x) = 0$ for all $x \in X$. By Lemma 3.2, there exists for any $\epsilon > 0$ some positive integer N_0 such that when $m, n \geq N_0$,

$$|f_m(x) - f_n(x)| \leq \epsilon \quad \text{for all } x \in X.$$

Since $\lim_{n \rightarrow \infty} f_n(x) = 0$, we let n goes to infinity in the above inequality to obtain that when $m \geq N_0$,

$$|f_m(x)| \leq \epsilon \quad \text{for all } x \in X.$$

In other words, $\|f_m\|_{L^\infty(X)} \leq \epsilon$ when $m \geq N_0$, implying $\lim_{n \rightarrow \infty} \|f_n\|_{L^\infty(X)} = 0$. □

By Proposition 3.1 and Lemmas 3.2 and 3.3, we conclude our construction of RKBS with the ℓ^1 norm in the following result.

Theorem 3.4. *Let K be a bounded function on $X \times X$ that satisfies (A3). Then \mathcal{B} having the form (1.3) and \mathcal{B}^\sharp are RKBS on X with the reproducing kernel K .*

We shall discuss in the rest of this section conditions on translation invariant $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{C}$ for which Admissibility Assumption (A3) holds. Specifically, such K are of the form

$$K(\mathbf{s}, \mathbf{t}) = \int_{\mathbb{R}^d} e^{-i(\mathbf{s}-\mathbf{t}) \cdot \boldsymbol{\xi}} \varphi(\boldsymbol{\xi}) d\boldsymbol{\xi}, \quad \mathbf{s}, \mathbf{t} \in \mathbb{R}^d, \quad (3.3)$$

where $\mathbf{s} \cdot \mathbf{t}$ stands for the standard inner product on \mathbb{R}^d , and $\varphi \in L^1(\mathbb{R}^d)$, the space of Lebesgue integrable functions on \mathbb{R}^d . One should not confuse $L^1(\mathbb{R}^d)$ with $\ell^1(\mathbb{R}^d)$. The latter one is defined with respect to the counting measure on \mathbb{R}^d while the first one is with respect to the Lebesgue measure. Note that K is bounded and continuous on $\mathbb{R}^d \times \mathbb{R}^d$. We give a sufficient condition for so defined a function K to satisfy (A3).

Proposition 3.5. *Let K be given by (3.3). If φ is nonzero almost everywhere on \mathbb{R}^d with respect to the Lebesgue measure then K satisfies (A3).*

Proof. Suppose that there exists $\mathbf{c} \in \ell^1(\mathbb{N})$ and pairwise distinct points $\mathbf{s}_j \in \mathbb{R}^d$, $j \in \mathbb{N}$ such that

$$\sum_{j=1}^{\infty} c_j K(\mathbf{s}_j, \mathbf{t}) = 0 \text{ for all } \mathbf{t} \in \mathbb{R}^d.$$

This equation can be reformulated by (3.3) as

$$\int_{\mathbb{R}^d} \left(\sum_{j=1}^{\infty} c_j e^{-i\mathbf{s}_j \cdot \boldsymbol{\xi}} \right) \varphi(\boldsymbol{\xi}) e^{i\mathbf{t} \cdot \boldsymbol{\xi}} d\boldsymbol{\xi} = 0 \text{ for all } \mathbf{t} \in \mathbb{R}^d.$$

It follows that for almost every $\boldsymbol{\xi} \in \mathbb{R}^d$ with respect to the Lebesgue measure

$$\left(\sum_{j=1}^{\infty} c_j e^{-i\mathbf{s}_j \cdot \boldsymbol{\xi}} \right) \varphi(\boldsymbol{\xi}) = 0.$$

By the assumption on φ ,

$$\sum_{j=1}^{\infty} c_j e^{-i\mathbf{s}_j \cdot \boldsymbol{\xi}} = 0 \text{ for almost every } \boldsymbol{\xi} \in \mathbb{R}^d.$$

Note that the function on the left hand side above is continuous on $\boldsymbol{\xi}$. We hence obtain that the Fourier transform of the discrete measure

$$\nu(A) := \sum_{\mathbf{s}_j \in A} c_j \text{ for every Borel subset } A \subseteq \mathbb{R}^d$$

is zero. Consequently, ν is the zero measure, implying $\mathbf{c} = 0$. \square

We next present a particular example as a corollary to Proposition 3.5.

Corollary 3.6. *If ϕ is nontrivial continuous function on \mathbb{R}^d with a compact support then $K(\mathbf{s}, \mathbf{t}) = \phi(\mathbf{s} - \mathbf{t})$, $\mathbf{s}, \mathbf{t} \in \mathbb{R}^d$ satisfies (A3).*

Proof. We regard ϕ as a tempered distribution and note by the Paley-Wiener theorem that the Fourier transform of ϕ is real-analytic on \mathbb{R}^d . Therefore, the Fourier transform of ϕ is nonzero everywhere on \mathbb{R}^d except at a subset of zero Lebesgue measure. The arguments similar to those in the proof of the last proposition hence apply. \square

We next present by Proposition 3.5 and Corollary 3.6 several examples of K that satisfy (A3) and hence can be used to construct RKBS with the ℓ^1 norm. Such functions include:

- the exponential kernel

$$K(\mathbf{s}, \mathbf{t}) = \exp(-\|\mathbf{s} - \mathbf{t}\|_{\ell^1(\mathbb{N}_d)}) = \frac{1}{\pi^d} \int_{\mathbb{R}^d} e^{-i(\mathbf{s}-\mathbf{t}) \cdot \boldsymbol{\xi}} \prod_{j=1}^d \frac{1}{1 + \xi_j^2} d\boldsymbol{\xi}, \quad \mathbf{s}, \mathbf{t} \in \mathbb{R}^d,$$

where for $\mathbf{s} \in \mathbb{R}^d$, $\|\mathbf{s}\|_2$ is its standard Euclidean norm on \mathbb{R}^d .

- the Gaussian kernel

$$K(\mathbf{s}, \mathbf{t}) = \exp\left(-\frac{\|\mathbf{s} - \mathbf{t}\|_2^2}{\sigma}\right) = \left(\frac{\sqrt{\sigma}}{2\sqrt{\pi}}\right)^d \int_{\mathbb{R}^d} e^{-i(\mathbf{s}-\mathbf{t}) \cdot \boldsymbol{\xi}} \exp\left(-\frac{\sigma}{4}\|\boldsymbol{\xi}\|_2^2\right) d\boldsymbol{\xi}, \quad \mathbf{s}, \mathbf{t} \in \mathbb{R}^d. \quad (3.4)$$

- inverse multiquadrics

$$K(\mathbf{s}, \mathbf{t}) = \left(\frac{1}{1 + \|\mathbf{s} - \mathbf{t}\|_2^2}\right)^\beta, \quad \mathbf{s}, \mathbf{t} \in \mathbb{R}^d, \quad \beta > 0, \quad (3.5)$$

whose Fourier transform is given by the modified Bessel function and is positive almost everywhere on \mathbb{R}^d (see [28], pages 52, 76 and 95).

- B-spline kernels

$$K(\mathbf{s}, \mathbf{t}) = \prod_{j=1}^d B_p(s_j - t_j), \quad \mathbf{s}, \mathbf{t} \in \mathbb{R}^d,$$

where s_j is the j -th component of \mathbf{s} and B_p denotes the p -th order B-spline, $p \geq 2$. B-spline kernels satisfies (A3) as they are given by bounded continuous functions of compact support.

- radial basis functions of compact support, including Wu's functions [29] and Wendland's functions [28]. Such functions are of the form $K(\mathbf{s}, \mathbf{t}) = \phi(\|\mathbf{s} - \mathbf{t}\|_2)$, $\mathbf{s}, \mathbf{t} \in \mathbb{R}^d$, where ϕ is a compactly supported univariate function dependent on the dimension d . We give two examples for $d = 3$:

$$\phi(r) := (1 - r)_+^2 \text{ and } \phi(r) := (1 - r)_+^4(1 + 4r), \quad r \geq 0$$

where $t_+ := \max\{0, t\}$ for $t \in \mathbb{R}$. These functions satisfy (A3) by Corollary 3.6.

On the other hand, a translation invariant K does not satisfy (A3) if its Fourier transform is compactly supported, as indicated in the next result.

Proposition 3.7. *If $\varphi \in L^1(\mathbb{R}^d)$ is compactly supported on \mathbb{R}^d then K given by (3.3) does not satisfy (A3).*

Proof. Without loss of generality, we may assume that $\text{supp } \varphi \subseteq [-1, 1]^d$. Choose a nontrivial infinitely continuously differentiable function ϕ that is supported on $[-\pi, \pi]^d$ and vanishes on $[-1, 1]^d$. We expand ϕ to a Fourier series

$$\phi(\boldsymbol{\xi}) = \sum_{\mathbf{j} \in \mathbb{Z}^d} c_{\mathbf{j}} e^{-i\mathbf{j} \cdot \boldsymbol{\xi}}, \quad \boldsymbol{\xi} \in [-\pi, \pi]^d,$$

where $c_{\mathbf{j}}$ is the Fourier coefficient of ϕ . Note that $\{c_{\mathbf{j}} : \mathbf{j} \in \mathbb{Z}^d\} \in \ell^1(\mathbb{Z}^d)$ as ϕ is infinitely continuously differentiable on $[-\pi, \pi]^d$. By arguments in the proof of Proposition 3.5,

$$\sum_{\mathbf{j} \in \mathbb{Z}^d} c_{\mathbf{j}} K(\mathbf{j}, \mathbf{t}) = \int_{\mathbb{R}^d} \left(\sum_{\mathbf{j} \in \mathbb{Z}^d} c_{\mathbf{j}} e^{-i\mathbf{j} \cdot \boldsymbol{\xi}} \right) \varphi(\boldsymbol{\xi}) e^{i\mathbf{t} \cdot \boldsymbol{\xi}} d\boldsymbol{\xi}, \quad \mathbf{t} \in \mathbb{R}^d.$$

By our construction,

$$\left(\sum_{\mathbf{j} \in \mathbb{Z}^d} c_{\mathbf{j}} e^{-i\mathbf{j} \cdot \boldsymbol{\xi}} \right) \varphi(\boldsymbol{\xi}) = 0 \text{ for all } \boldsymbol{\xi} \in \mathbb{R}^d,$$

which implies $\sum_{\mathbf{j} \in \mathbb{Z}^d} c_{\mathbf{j}} K(\mathbf{j}, \cdot) = 0$. Moreover, $c_{\mathbf{j}} \neq 0$ for at least one $\mathbf{j} \in \mathbb{Z}^d$ because ϕ is nontrivial. We obtain that K does not satisfy (A3). \square

By Proposition 3.7, the sinc kernel

$$K(\mathbf{s}, \mathbf{t}) := \text{sinc}(\mathbf{s} - \mathbf{t}) := \prod_{j=1}^d \frac{\sin(\pi(s_j - t_j))}{\pi(s_j - t_j)}, \quad \mathbf{s}, \mathbf{t} \in \mathbb{R}^d$$

does not satisfy (A3). As a consequence, it can not yield an RKBS with the ℓ^1 norm by the procedures introduced in this section. Similar arguments as those in the proof of Proposition 3.7 are able to show that if ν is a compactly supported Borel measure on \mathbb{R}^d of finite total variation then the following function

$$K(\mathbf{s}, \mathbf{t}) := \int_{\mathbb{R}^d} e^{-i(\mathbf{s}-\mathbf{t}) \cdot \boldsymbol{\xi}} d\nu(\boldsymbol{\xi}), \quad \mathbf{s}, \mathbf{t} \in \mathbb{R}^d$$

does not satisfy (A3). Instances include the class of Bessel-based radial functions [10] where the Borel measure is the dirac delta measure on the unit sphere of the Euclidean space.

4 Representer Theorems in RKBS with the ℓ^1 Norm

Up to now our arguments have relied on Admissibility Assumptions (A1)–(A3). In this section the final assumption, (A4), is invoked to guarantee that the representer theorem should hold for the constructed RKBS. A regularized learning scheme in the RKBS \mathcal{B} constructed by (1.3) can be generally expressed as finding f_0 such that

$$f_0 = \underset{f \in \mathcal{B}}{\text{argmin}} [V(f(\mathbf{x})) + \mu\phi(\|f\|_{\mathcal{B}})], \quad (4.1)$$

where $\mathbf{x} := \{x_j \in X : j \in \mathbb{N}_n\}$, $n \in \mathbb{N}$, is the sequence of given pairwise distinct sampling points, $f(\mathbf{x}) := (f(x_j) : j \in \mathbb{N}_n) \in \mathbb{C}^n$, $V : \mathbb{C}^n \rightarrow \mathbb{R}_+$ is a loss function, μ is a positive regularization parameter, and $\phi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ is a nondecreasing regularization function. Here, $\mathbb{R}_+ := [0, +\infty)$. The loss function and regularization function should satisfy some minimal requirements for the learning scheme (4.1) to be useful. This consideration gives rise to the following definition.

Definition 4.1. *A regularized learning scheme (4.1) is said to be acceptable if V and ϕ are continuous and*

$$\lim_{t \rightarrow \infty} \phi(t) = +\infty. \quad (4.2)$$

It is possible that the solution to (4.1) is non-unique, and in that case we are only interested in finding one possible solution.

We now introduce the main concept of this section.

Definition 4.2. *The space \mathcal{B} is said to satisfy the linear representer theorem for regularized learning if every acceptable regularized learning scheme (4.1) has a minimizer of the form*

$$f_0 = \sum_{j=1}^n c_j K(x_j, \cdot), \quad (4.3)$$

where c_j 's are constants. In other words, there exists a solution f_0 lying in the finite dimensional subspace $\mathcal{S}^{\mathbf{x}} := \text{span}\{K(x_j, \cdot) : j \in \mathbb{N}_n\}$.

An RKHS with K being its reproducing kernel in the usual sense always satisfies the linear representer theorem [14]. The result for uniformly convex and uniformly Fréchet differentiable pre-RKBS with a reproducing kernel given by the semi-inner product was established in [31, 32]. For more information on this important property for RKHS and vector-valued RKHS, see, for example, [1, 17, 21] and the references cited therein.

Our purpose is to discuss the conditions on K such that \mathcal{B} satisfies the linear representer theorem. The representer theorem for (4.1) is closely related to the representer theorem for the minimal norm interpolation problem. In the RKHS case, an equivalence was proved in [16]. We shall follow the approach to consider the minimal norm interpolation in \mathcal{B} first. For any $\mathbf{y} \in \mathbb{C}^n$, set $\mathcal{I}_{\mathbf{x}}(\mathbf{y})$ to be the subset of functions in \mathcal{B} that interpolate the specified data, namely, $\mathcal{I}_{\mathbf{x}}(\mathbf{y}) := \{f \in \mathcal{B} : f(\mathbf{x}) = \mathbf{y}\}$. A minimal norm interpolant in \mathcal{B} is a function f_{\min} satisfying

$$f_{\min} = \operatorname{argmin}\{\|f\|_{\mathcal{B}} : f \in \mathcal{I}_{\mathbf{x}}(\mathbf{y})\}. \quad (4.4)$$

Again, in the case of a non-unique solution, we are only interested in obtaining one solution. Since $K[\mathbf{x}]$ is nonsingular, one sees that the typically infinite dimensional $\mathcal{I}_{\mathbf{x}}(\mathbf{y})$ always has a non-empty intersection with $\mathcal{S}^{\mathbf{x}}$, for all $\mathbf{y} \in \mathbb{C}^n$ and pairwise distinct $\mathbf{x} \subseteq X$.

Definition 4.3. *An RKBS \mathcal{B} is said to satisfy the linear representer theorem for minimal norm interpolation if for any choice of data, \mathbf{x} and \mathbf{y} , there is a minimal norm interpolant, (4.4), lying in $\mathcal{S}^{\mathbf{x}}$.*

We shall show that \mathcal{B} satisfies the linear representer theorem for regularized learning if and only if it does so for minimal norm interpolation. We first prove one direction of the equivalence.

Lemma 4.4. *If \mathcal{B} satisfies the linear representer theorem for the minimal norm interpolation, then it also does so for regularized learning.*

Proof. Let V , ϕ , and μ be arbitrary, but fixed according to the conditions that (4.1) be an acceptable regularization scheme. For an arbitrary function f in \mathcal{B} . We let f_0 be the minimizer of $\inf_{g \in \mathcal{I}_{\mathbf{x}}(f(\mathbf{x}))} \|g\|_{\mathcal{B}}$ that has the form (4.3). Then $f_0(\mathbf{x}) = f(\mathbf{x})$ and $\|f_0\|_{\mathcal{B}} \leq \|f\|_{\mathcal{B}}$. As a consequence, $V(f_0(\mathbf{x})) = V(f(\mathbf{x}))$ but $\phi(\|f_0\|_{\mathcal{B}}) \leq \phi(\|f\|_{\mathcal{B}})$ as ϕ is nondecreasing. It follows that

$$\inf_{f \in \mathcal{B}} V(f(\mathbf{x})) + \mu\phi(\|f\|_{\mathcal{B}}) = \inf_{f \in \mathcal{S}^{\mathbf{x}}} V(f(\mathbf{x})) + \mu\phi(\|f\|_{\mathcal{B}}).$$

By (4.2), there exists a positive constant α such that

$$\inf_{f \in \mathcal{S}^{\mathbf{x}}} V(f(\mathbf{x})) + \mu\phi(\|f\|_{\mathcal{B}}) = \inf_{f \in \mathcal{S}^{\mathbf{x}}, \|f\|_{\mathcal{B}} \leq \alpha} V(f(\mathbf{x})) + \mu\phi(\|f\|_{\mathcal{B}}).$$

Note that the functional we are minimizing is continuous on \mathcal{B} by the assumption on V , ϕ and by the continuity of point evaluation functionals on \mathcal{B} . By the elementary fact that a continuous function on a compact metric space attains its minimum in the space, (4.1) has a minimizer that belongs to $\{f \in \mathcal{S}^{\mathbf{x}} : \|f\|_{\mathcal{B}} \leq \alpha\}$. Therefore, \mathcal{B} satisfies the linear representer theorem. \square

For the other direction, it suffices to consider a class of regularization functionals with a particular choice of V and ϕ . In the limit of vanishing μ we recover the minimal norm interpolant.

Lemma 4.5. *If \mathcal{B} satisfies the linear representer theorem for regularized learning, then it also satisfies the linear representer theorem for minimal norm interpolation.*

Proof. We shall follow the idea in [16]. Choose any $n \in \mathbb{N}_n$, any $\mathbf{x} = \{x_j \in X : j \in \mathbb{N}_n\}$ with pairwise distinct elements, and any $\mathbf{y} \in \mathbb{C}^n$. For every $\mu > 0$, let $f_{0,\mu} \in \mathcal{S}^{\mathbf{x}}$ be a minimizer of (4.1) with the choice of

$$V(f(\mathbf{x})) = \|f(\mathbf{x}) - \mathbf{y}\|_2^2, \quad \phi(t) = t. \quad (4.5)$$

Here, $\|\cdot\|_2$ is the standard Euclidean norm on \mathbb{C}^n . Defining the $1 \times n$ row vector function by

$$K^{\mathbf{x}}(x) := (K(x_j, x) : j \in \mathbb{N}_n) \text{ for all } x \in X.$$

It follows that $f_{0,\mu} = K^{\mathbf{x}}(\cdot)\mathbf{c}_\mu$ for some $\mathbf{c}_\mu \in \mathbb{C}^n$. Then we have

$$\|K[\mathbf{x}]\mathbf{c}_\mu - \mathbf{y}\|_2^2 = \|f_{0,\mu}(\mathbf{x}) - \mathbf{y}\|_2^2 \leq V(f_{0,\mu}) + \mu\phi(\|f_{0,\mu}\|_{\mathcal{B}}) \leq V(0) + \mu\phi(\|0\|_{\mathcal{B}}) = \|\mathbf{y}\|_2^2.$$

As $K[\mathbf{x}]$ is nonsingular, the above inequality implies that $\{\mathbf{c}_\mu : \mu > 0\}$ forms a bounded set in \mathbb{C}^n . By restricting to a subsequence if necessary, we may hence assume that \mathbf{c}_μ converges to some $\mathbf{c}_0 \in \mathbb{C}^n$ as μ goes to zero. We shall show that $f_{0,0} := K^{\mathbf{x}}(\cdot)\mathbf{c}_0 \in \mathcal{S}^{\mathbf{x}}$ is a minimal norm interpolant.

Since \mathbf{c}_μ converges to \mathbf{c}_0 as μ tends to zero, we first get

$$\lim_{\mu \rightarrow 0} \|f_{0,\mu} - f_{0,0}\|_{\mathcal{B}} = \lim_{\mu \rightarrow 0} \|\mathbf{c}_\mu - \mathbf{c}_0\|_{\ell^1(\mathbb{N}_n)} = 0. \quad (4.6)$$

Since point evaluation functionals are continuous on \mathcal{B} , we obtain by (4.6)

$$f_{0,0}(x_j) = \lim_{\mu \rightarrow 0} f_{0,\mu}(x_j) \text{ for all } j \in \mathbb{N}_n. \quad (4.7)$$

Now let g be an arbitrary interpolant, i.e., an arbitrary element of $\mathcal{I}_{\mathbf{x}}(\mathbf{y})$. As $f_{0,\mu}$ is a minimizer of (4.1) with the choice (4.5), it follows that

$$\|f_{0,\mu}(\mathbf{x}) - \mathbf{y}\|_2^2 + \mu\|f_{0,\mu}\|_{\mathcal{B}} \leq \|g(\mathbf{x}) - \mathbf{y}\|_2^2 + \mu\|g\|_{\mathcal{B}} = \mu\|g\|_{\mathcal{B}}. \quad (4.8)$$

Letting $\mu \rightarrow 0$ on both sides of the above inequality, we obtain by (4.7) $\|f_{0,0}(\mathbf{x}) - \mathbf{y}\|_2^2 = 0$, which implies that $f_{0,0}$ is also an interpolant, i.e., $f_{0,0} \in \mathcal{I}_{\mathbf{x}}(\mathbf{y})$. It also follows from (4.8) that $\|f_{0,\mu}\|_{\mathcal{B}} \leq \|g\|_{\mathcal{B}}$ for all $\mu > 0$, which together with (4.6) implies $\|f_{0,0}\|_{\mathcal{B}} \leq \|g\|_{\mathcal{B}}$. Since g is an arbitrary function in $\mathcal{I}_{\mathbf{x}}(\mathbf{y})$ and $f_{0,0} \in \mathcal{I}_{\mathbf{x}}(\mathbf{y})$, we see that $f_{0,0}$ is a minimal norm interpolant, i.e., a solution of (4.4). The proof is complete. \square

Combining Lemmas 4.4 and 4.5, we reach the characterization for \mathcal{B} to satisfy the linear representer theorem.

Proposition 4.6. *The space \mathcal{B} satisfies the linear representer theorem for regularized learning if and only if \mathcal{B} satisfies the linear representer theorem for minimal norm interpolation.*

In view of the above result, we shall focus on necessary and sufficient conditions for the minimal norm interpolation in \mathcal{B} to satisfy the linear representer theorem. To this end, we begin with the simplest case when only one more sampling point is added to \mathbf{x} . Recall the definition of $K_{\mathbf{x}}(x)$ from the introduction. It is worthwhile to point out that $K_{\mathbf{x}}(x)$ is in general not the transpose of $K^{\mathbf{x}}(x)$ as K is not required to be symmetric.

Lemma 4.7. *Let $\mathbf{x} = \{x_j \in X : j \in \mathbb{N}_n\}$ have pairwise distinct elements, let x_{n+1} be an arbitrary point in $X \setminus \mathbf{x}$, and set $\bar{\mathbf{x}} := \{x_j : j \in \mathbb{N}_{n+1}\}$. It follows that the minimum norm interpolant in $\mathcal{S}^{\bar{\mathbf{x}}}$ is the same as the minimum norm interpolant in $\mathcal{S}^{\mathbf{x}}$, i.e.,*

$$\min_{f \in \mathcal{I}_{\mathbf{x}}(\mathbf{y}) \cap \mathcal{S}^{\bar{\mathbf{x}}}} \|f\|_{\mathcal{B}} = \min_{f \in \mathcal{I}_{\mathbf{x}}(\mathbf{y}) \cap \mathcal{S}^{\mathbf{x}}} \|f\|_{\mathcal{B}} \text{ for all } \mathbf{y} \in \mathbb{C}^n, \quad (4.9)$$

if and only if (1.2) holds true.

Proof. Notice that $\mathcal{I}_{\mathbf{x}}(\mathbf{y}) \cap \mathcal{S}^{\mathbf{x}}$ has only one function $f = K^{\mathbf{x}}(\cdot)K[\mathbf{x}]^{-1}\mathbf{y}$. We next estimate the norm of functions in $\mathcal{I}_{\mathbf{x}}(\mathbf{y}) \cap \mathcal{S}^{\overline{\mathbf{x}}}$. Let $g \in \mathcal{I}_{\mathbf{x}}(\mathbf{y}) \cap \mathcal{S}^{\overline{\mathbf{x}}}$ and $b := g(x_{n+1})$. Note that g is uniquely determined by b as it has already satisfied the interpolation condition $g(\mathbf{x}) = \mathbf{y}$. In fact, as $K[\overline{\mathbf{x}}]$ is nonsingular, $g = K^{\overline{\mathbf{x}}}(\cdot)K[\overline{\mathbf{x}}]^{-1}\overline{\mathbf{y}}$, where $\overline{\mathbf{y}} = (\mathbf{y}^T, b)^T \in \mathbb{C}^{n+1}$. Direct computations show that

$$K[\overline{\mathbf{x}}]^{-1}\overline{\mathbf{y}} = \begin{pmatrix} K[\mathbf{x}] & K_{\mathbf{x}}(x_{n+1}) \\ K^{\mathbf{x}}(x_{n+1}) & K(x_{n+1}, x_{n+1}) \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{y} \\ b \end{pmatrix} = \begin{pmatrix} K[\mathbf{x}]^{-1}\mathbf{y} + \frac{q}{p}K[\mathbf{x}]^{-1}K_{\mathbf{x}}(x_{n+1}) \\ -\frac{q}{p} \end{pmatrix},$$

where $p := K(x_{n+1}, x_{n+1}) - K^{\mathbf{x}}(x_{n+1})K[\mathbf{x}]^{-1}K_{\mathbf{x}}(x_{n+1})$ and $q := K^{\mathbf{x}}(x_{n+1})K[\mathbf{x}]^{-1}\mathbf{y} - b$.

We now show sufficiency. If (1.2) holds true then we have

$$\begin{aligned} \|g\|_{\mathcal{B}} &= \|K[\overline{\mathbf{x}}]^{-1}\overline{\mathbf{y}}\|_{\ell^1(\mathbb{N}_{n+1})} \geq \|K[\mathbf{x}]^{-1}\mathbf{y}\|_{\ell^1(\mathbb{N}_n)} - \|(K[\mathbf{x}])^{-1}K_{\mathbf{x}}(x_{n+1})\|_{\ell^1(\mathbb{N}_n)} \left|\frac{q}{p}\right| + \left|\frac{q}{p}\right| \\ &\geq \|K[\mathbf{x}]^{-1}\mathbf{y}\|_{\ell^1(\mathbb{N}_n)} = \|f\|_{\mathcal{B}}, \end{aligned}$$

which implies

$$\min_{f \in \mathcal{I}_{\mathbf{x}}(\mathbf{y}) \cap \mathcal{S}^{\overline{\mathbf{x}}}} \|f\|_{\mathcal{B}} \geq \min_{f \in \mathcal{I}_{\mathbf{x}}(\mathbf{y}) \cap \mathcal{S}^{\mathbf{x}}} \|f\|_{\mathcal{B}}.$$

Since $\mathcal{S}^{\mathbf{x}} \subseteq \mathcal{S}^{\overline{\mathbf{x}}}$,

$$\min_{f \in \mathcal{I}_{\mathbf{x}}(\mathbf{y}) \cap \mathcal{S}^{\overline{\mathbf{x}}}} \|f\|_{\mathcal{B}} \leq \min_{f \in \mathcal{I}_{\mathbf{x}}(\mathbf{y}) \cap \mathcal{S}^{\mathbf{x}}} \|f\|_{\mathcal{B}}.$$

Thus, (4.9) holds true.

On the other hand, if (4.9) is always true for all $\mathbf{y} \in \mathbb{C}^n$ then we must have

$$\|K[\overline{\mathbf{x}}]^{-1}\overline{\mathbf{y}}\|_{\ell^1(\mathbb{N}_{n+1})} \geq \|K[\mathbf{x}]^{-1}\mathbf{y}\|_{\ell^1(\mathbb{N}_n)} \text{ for all } \mathbf{y} \in \mathbb{C}^n \text{ and } b \in \mathbb{C}.$$

In particular, the choices $\mathbf{y} = K_{\mathbf{x}}(x_{n+1})$ and $b = K^{\mathbf{x}}(x_{n+1})K[\mathbf{x}]^{-1}K_x^T(x_{n+1}) + p$ yields that

$$\|K[\overline{\mathbf{x}}]^{-1}\overline{\mathbf{y}}\|_{\ell^1(\mathbb{N}_{n+1})} = \left\| \begin{pmatrix} \mathbf{0} \\ 1 \end{pmatrix} \right\|_{\ell^1(\mathbb{N}_{n+1})} = 1 \quad \text{and} \quad \|K[\mathbf{x}]^{-1}\mathbf{y}\|_{\ell^1(\mathbb{N}_n)} = \|(K[\mathbf{x}])^{-1}K_{\mathbf{x}}(x_{n+1})\|_{\ell^1(\mathbb{N}_n)}.$$

Combing the above two equations proves (1.2). The proof is complete. \square

We are now ready to present one of the main results in this paper.

Theorem 4.8. *Every minimal norm interpolant (4.4) in \mathcal{B} satisfies the linear representer theorem if and only if (1.2) holds true for all $n \in \mathbb{N}$ and all pairwise distinct sampling points $x_j \in X$, $j \in \mathbb{N}_{n+1}$.*

Proof. The minimal norm interpolant (4.4) satisfies the linear representer theorem if and only if

$$\min_{g \in \mathcal{I}_{\mathbf{x}}(\mathbf{y})} \|g\|_{\mathcal{B}} = \min_{f \in \mathcal{I}_{\mathbf{x}}(\mathbf{y}) \cap \mathcal{S}^{\mathbf{x}}} \|f\|_{\mathcal{B}}.$$

Therefore, if the above equation holds true then since $\mathcal{I}_{\mathbf{x}}(\mathbf{y}) \cap \mathcal{S}^{\mathbf{x}} \subseteq \mathcal{I}_{\mathbf{x}}(\mathbf{y}) \cap \mathcal{S}^{\overline{\mathbf{x}}} \subseteq \mathcal{I}_{\mathbf{x}}(\mathbf{y})$, we obtain (4.9). By Lemma 4.7, (1.2) is true for every $x_{n+1} \in X$.

It remains to prove the sufficiency. We shall first show $\|g\|_{\mathcal{B}} \geq \min_{f \in \mathcal{I}_{\mathbf{x}}(\mathbf{y}) \cap \mathcal{S}^{\mathbf{x}}} \|f\|_{\mathcal{B}}$ for all $g \in \mathcal{I}_{\mathbf{x}}(\mathbf{y}) \cap \mathcal{B}_0$. To this end, we express g as $g = \sum_{j=1}^m c_j K(x_j, \cdot)$ for some $m \geq n$ and pairwise distinct $\{x_j : j \in \mathbb{N}_m\} \subseteq X$. This can always be done by adding some sampling points, setting the corresponding coefficients to be zero, and relabeling if necessary. We let $y_j := g(x_j)$, $j \in \mathbb{N}_m$,

$\mathbf{u}_l := (y_j : j \in \mathbb{N}_l)$, and $\mathbf{v}_l = \{x_j : j \in \mathbb{N}_l\}$ for $1 \leq l \leq m$. Note that $\mathbf{y} = \mathbf{u}_n$ and $\mathbf{x} = \mathbf{v}_n$. It follows that $g \in \mathcal{I}_{\mathbf{v}_m}(\mathbf{u}_m) \cap \mathcal{S}^{\mathbf{v}_m}$ and thus,

$$\|g\|_{\mathcal{B}} \geq \min_{f \in \mathcal{I}_{\mathbf{v}_m}(\mathbf{u}_m) \cap \mathcal{S}^{\mathbf{v}_m}} \|f\|_{\mathcal{B}}.$$

Since $\mathcal{I}_{\mathbf{v}_m}(\mathbf{u}_m) \subseteq \mathcal{I}_{\mathbf{v}_{m-1}}(\mathbf{u}_{m-1})$, we apply Lemma 4.7 to get

$$\min_{f \in \mathcal{I}_{\mathbf{v}_m}(\mathbf{u}_m) \cap \mathcal{S}^{\mathbf{v}_m}} \|f\|_{\mathcal{B}} \geq \min_{f \in \mathcal{I}_{\mathbf{v}_{m-1}}(\mathbf{u}_{m-1}) \cap \mathcal{S}^{\mathbf{v}_m}} \|f\|_{\mathcal{B}} = \min_{f \in \mathcal{I}_{\mathbf{v}_{m-1}}(\mathbf{u}_{m-1}) \cap \mathcal{S}^{\mathbf{v}_{m-1}}} \|f\|_{\mathcal{B}}.$$

It follows that

$$\|g\|_{\mathcal{B}} \geq \min_{f \in \mathcal{I}_{\mathbf{v}_{m-1}}(\mathbf{u}_{m-1}) \cap \mathcal{S}^{\mathbf{v}_{m-1}}} \|f\|_{\mathcal{B}}.$$

Repeating this process, we reach

$$\|g\|_{\mathcal{B}} \geq \min_{f \in \mathcal{I}_{\mathbf{v}_n}(\mathbf{u}_n) \cap \mathcal{S}^{\mathbf{v}_n}} \|f\|_{\mathcal{B}} = \min_{f \in \mathcal{I}_{\mathbf{x}}(\mathbf{y}) \cap \mathcal{S}^{\mathbf{x}}} \|f\|_{\mathcal{B}} \text{ for all } g \in \mathcal{I}_{\mathbf{x}}(\mathbf{y}) \cap \mathcal{B}_0. \quad (4.10)$$

Now let $g \in \mathcal{I}_{\mathbf{x}}(\mathbf{y})$ be arbitrary but fixed. Then there exists a sequence of functions $\{g_j \in \mathcal{B}_0 : j \in \mathbb{N}\}$ that converges to g in \mathcal{B} . We let f and f_j be the function in $\mathcal{S}^{\mathbf{x}}$ such that $f(\mathbf{x}) = \mathbf{y}$ and $f_j(\mathbf{x}) = g_j(\mathbf{x})$, $j \in \mathbb{N}$. They are explicitly given by

$$f = K^{\mathbf{x}}(\cdot)K[\mathbf{x}]^{-1}g(\mathbf{x}) \quad \text{and} \quad f_j = K^{\mathbf{x}}(\cdot)K[\mathbf{x}]^{-1}g_j(\mathbf{x}), \quad j \in \mathbb{N}.$$

Since g_j converges to g in \mathcal{B} and point evaluation functionals are continuous on \mathcal{B} , $g_j(\mathbf{x}) \rightarrow g(\mathbf{x})$ as $j \rightarrow \infty$. As a result, $\lim_{j \rightarrow \infty} \|f - f_j\|_{\mathcal{B}} = 0$. By (4.10), $\|g_j\|_{\mathcal{B}} \geq \|f_j\|_{\mathcal{B}}$ for all $j \in \mathbb{N}$. We hence obtain that $\|g\|_{\mathcal{B}} \geq \|f\|_{\mathcal{B}}$. Therefore,

$$\min_{g \in \mathcal{I}_{\mathbf{x}}(\mathbf{y})} \|g\|_{\mathcal{B}} \geq \min_{f \in \mathcal{I}_{\mathbf{x}}(\mathbf{y}) \cap \mathcal{S}^{\mathbf{x}}} \|f\|_{\mathcal{B}}.$$

The reverse direction of the inequality is clear as $\mathcal{I}_{\mathbf{x}}(\mathbf{y}) \cap \mathcal{S}^{\mathbf{x}} \subseteq \mathcal{I}_{\mathbf{x}}(\mathbf{y})$. □

We draw the following conclusion by Theorems 4.6 and 4.8.

Corollary 4.9. *Every acceptable regularized learning scheme of the form (4.1) has a minimizer of the form (4.3) if and only if the function K satisfies the property (1.2).*

In the last part of the section, we briefly discuss the linear representer theorem in \mathcal{B}^{\sharp} under the same assumption that K is bounded and satisfies (A3). By Theorem 3.4, \mathcal{B}^{\sharp} is an RKBS on X . Likewise, we call a regularized learning scheme

$$f_0 = \operatorname{argmin}_{f \in \mathcal{B}^{\sharp}} V(f(\mathbf{x})) + \mu\phi(\|f\|_{\mathcal{B}^{\sharp}}) \quad (4.11)$$

acceptable if V and ϕ are continuous and (4.2) is satisfied by ϕ . The space \mathcal{B}^{\sharp} is said to satisfy the linear representer theorem if every acceptable learning scheme (4.11) has a minimizer of the following form

$$f_0 = \sum_{j=1}^n c_j K(\cdot, x_j), \quad (4.12)$$

where c_j 's are constants. We follow similar approaches to those used for \mathcal{B} to study this important property on \mathcal{B}^{\sharp} .

Proposition 4.10. *Let $\mathbf{x} \subseteq X$ have pairwise distinct elements. Every acceptable regularized learning scheme (4.11) in \mathcal{B}^\sharp has a minimizer, f_0 lying in $\mathcal{S}_{\mathbf{x}} := \text{span}\{K(\cdot, x_j) : j \in \mathbb{N}_n\}$ if and only if there is a minimal norm interpolant,*

$$f_{\min} := \underset{f \in \mathcal{B}^\sharp, f(\mathbf{x}) = \mathbf{y}}{\operatorname{argmin}} \|f\|_{\mathcal{B}^\sharp} \quad (4.13)$$

lying in $\mathcal{S}_{\mathbf{x}}$ for all $\mathbf{y} \in \mathbb{C}^n$.

Proof. The arguments of the proof are similar to those for \mathcal{B} . One only needs to note that although the norm of a function in \mathcal{B}^\sharp may not be known, any two norms on the finite dimensional vector space $\mathcal{S}_{\mathbf{x}}$ are equivalent. \square

To study conditions ensuring that the minimal norm interpolation (4.13) satisfies the linear representer theorem, we first identify a specific form of the norm $\|\cdot\|_{\mathcal{B}^\sharp}$ under the assumption that K satisfies (1.2). Notice that a function $f_{\mathbf{c}} = \sum_{j=1}^n c_j K(\cdot, x_j) \in \mathcal{S}_{\mathbf{x}} \subseteq \mathcal{B}_0^\sharp$ can be represented as $f_{\mathbf{c}} = \mathbf{c}^T K_{\mathbf{x}}(\cdot)$.

Lemma 4.11. *Let \mathbf{x} have pairwise distinct elements. The function K satisfies (1.2) if and only if*

$$\|f_{\mathbf{c}}\|_{\mathcal{B}^\sharp} = \|\mathbf{c}^T K[\mathbf{x}]\|_\infty \text{ for all } f_{\mathbf{c}} = \mathbf{c}^T K_{\mathbf{x}}(\cdot), \quad \mathbf{c} \in \mathbb{C}^n, \quad (4.14)$$

where $\|\cdot\|_\infty$ denotes the maximum norm on \mathbb{C}^n .

Proof. Suppose that K satisfies (1.2) for all $x_{n+1} \in X \setminus \mathbf{x}$. Then we have for all $x \in X$ that $\|K[\mathbf{x}]^{-1} K_{\mathbf{x}}(x)\|_{\ell^1(\mathbb{N}_n)} \leq 1$. Let $\mathbf{c} \in \mathbb{C}^n$ and $x \in X$. It follows from this inequality that

$$|\mathbf{c}^T K_{\mathbf{x}}(x)| = |\mathbf{c}^T K[\mathbf{x}] K[\mathbf{x}]^{-1} K_{\mathbf{x}}(x)| \leq \|\mathbf{c}^T K[\mathbf{x}]\|_\infty \|K[\mathbf{x}]^{-1} K_{\mathbf{x}}(x)\|_{\ell^1(\mathbb{N}_n)} \leq \|\mathbf{c}^T K[\mathbf{x}]\|_\infty,$$

which implies by Lemma 3.2 that for $f_{\mathbf{c}} = \mathbf{c}^T K_{\mathbf{x}}(\cdot)$

$$\|f_{\mathbf{c}}\|_{\mathcal{B}^\sharp} = \|\mathbf{c}^T K_{\mathbf{x}}(\cdot)\|_{L^\infty(X)} \leq \|\mathbf{c}^T K[\mathbf{x}]\|_\infty.$$

The other direction of the inequality is clear as we have

$$\|\mathbf{c}^T K[\mathbf{x}]\|_\infty = \max\{|\mathbf{c}^T K_{\mathbf{x}}(x_j)| : j \in \mathbb{N}_n\} \leq \|\mathbf{c}^T K_{\mathbf{x}}(\cdot)\|_{L^\infty(X)} = \|f_{\mathbf{c}}\|_{\mathcal{B}^\sharp}.$$

It remains to show that (4.14) implies (1.2). We prove this by construction. For any $x_{n+1} \in X$, we can find a nonzero vector $\mathbf{c} \in \mathbb{C}^n$ such that

$$|\mathbf{c}^T K_{\mathbf{x}}(x_{n+1})| = |\mathbf{c}^T K[\mathbf{x}] K[\mathbf{x}]^{-1} K_{\mathbf{x}}(x_{n+1})| = \|\mathbf{c}^T K[\mathbf{x}]\|_\infty \|K[\mathbf{x}]^{-1} K_{\mathbf{x}}(x_{n+1})\|_{\ell^1(\mathbb{N}_n)}.$$

We then let $f_{\mathbf{c}} = \mathbf{c}^T K_{\mathbf{x}}(\cdot)$ and obtain by (4.14)

$$\|\mathbf{c}^T K[\mathbf{x}]\|_\infty \|K[\mathbf{x}]^{-1} K_{\mathbf{x}}(x_{n+1})\|_{\ell^1(\mathbb{N}_n)} = |f_{\mathbf{c}}(x_{n+1})| \leq \|f_{\mathbf{c}}\|_{L^\infty(X)} = \|f_{\mathbf{c}}\|_{\mathcal{B}^\sharp} = \|\mathbf{c}^T K[\mathbf{x}]\|_\infty,$$

which implies (1.2) for $\mathbf{c}^T K[\mathbf{x}]$ is not the zero vector. The proof is complete. \square

We now show that (1.2) is sufficient for \mathcal{B}^\sharp to satisfy the linear representer theorem.

Theorem 4.12. *If K satisfies (1.2) then \mathcal{B}^\sharp satisfies the linear representer theorem.*

Proof. Suppose that (1.2) holds true. By Lemma 4.10, it suffices to show that the minimal norm interpolation (4.13) has a minimizer of the form (4.3). We shall prove this by directly showing that $f_0 = \mathbf{y}^T K[\mathbf{x}]^{-1} K_{\mathbf{x}}(\cdot)$ is a minimizer for (4.13). Let f be an arbitrary function in \mathcal{B}^\sharp such that $f(\mathbf{x}) = \mathbf{y}$. Then we have by Lemma 3.2

$$\|f\|_{\mathcal{B}^\sharp} = \|f\|_{L^\infty(X)} \geq \|f(\mathbf{x})\|_\infty = \|\mathbf{y}\|_\infty.$$

By Lemma 4.11,

$$\|f_0\|_{\mathcal{B}^\sharp} = \|\mathbf{y}^T K[\mathbf{x}]^{-1} K[\mathbf{x}]\|_\infty = \|\mathbf{y}\|_\infty.$$

Combining the above two inequalities leads to $\|f_0\|_{\mathcal{B}^\sharp} \leq \|f\|_{\mathcal{B}^\sharp}$. Therefore, (4.13) has the minimizer $f_0 = \mathbf{y}^T K[\mathbf{x}]^{-1} K_{\mathbf{x}}(\cdot)$ which has the form (4.12). \square

In the particular case when X has a finite cardinality, we shall show that condition (1.2) is also necessary for \mathcal{B}^\sharp to satisfy the linear representer theorem.

Proposition 4.13. *If X consists of finitely many points and \mathcal{B}^\sharp satisfies the linear representer theorem then (1.2) holds true.*

Proof. Let $\mathbf{c} \in \mathbb{C}^n$ and $f_{\mathbf{c}} = \mathbf{c}^T K_{\mathbf{x}}(\cdot)$. Under the assumptions, we get by Proposition 4.10 that $f_{\mathbf{c}}$ is a minimizer for the minimal norm interpolation (4.13) with $\mathbf{y} = f_{\mathbf{c}}(\mathbf{x}) = (K[\mathbf{x}])^T \mathbf{c}$. Since X has a finite cardinality and $K[\mathbf{x}]$ is nonsingular for all pairwise distinct $\mathbf{x} \subseteq X$, we can find a function $g \in \mathcal{B}_0$ such that $g(\mathbf{x}) = \mathbf{y}$ and $\|g\|_{L^\infty(X)} \leq \|\mathbf{y}\|_\infty$. Since $f_{\mathbf{c}}$ is a minimizer of (4.13) and g satisfies $g(\mathbf{x}) = \mathbf{y}$,

$$\|f_{\mathbf{c}}\|_{\mathcal{B}^\sharp} \leq \|g\|_{\mathcal{B}^\sharp} = \|g\|_{L^\infty(X)} = \|\mathbf{y}\|_\infty = \|(K[\mathbf{x}]^T) \mathbf{c}\|_\infty.$$

On the other hand, we have by Lemma 3.2

$$\|f_{\mathbf{c}}\|_{\mathcal{B}^\sharp} = \|f_{\mathbf{c}}\|_{L^\infty(X)} \geq \|f_{\mathbf{c}}(\mathbf{x})\|_\infty = \|(K[\mathbf{x}]^T) \mathbf{c}\|_\infty.$$

By the above two equations, (4.14) holds true. By Lemma 4.11, K satisfies (1.2). \square

One observes that the key ingredient in the proof of Proposition 4.13 is to extend a function on the discrete set \mathbf{x} to a function in \mathcal{B}^\sharp in a way that the supremum norm is preserved. In many cases, this is achievable without X being a finite set. For instance, by the Tietze extension theorem in topology, such an extension exists when X is a compact metric space and K is a universal kernel [19] on X . Thus, for those input spaces X and functions K , \mathcal{B}^\sharp satisfies the linear representer theorem if and only if (1.2) holds true.

5 Examples of Admissible Kernels

Recall the definition of admissible kernels from the introduction. Note that the first requirement (A1) in the definition implies (3.1). Theorem 1.2 is proved by combining Theorem 3.4 and Corollary 4.9. By this result, admissible kernels are crucial for our construction. Functions K satisfying requirements (A1)–(A3) are usually relatively easy to find. Some examples have been presented before Proposition 3.7 in Section 3. However, requirement (A4) could be somewhat demanding and rule out many commonly used kernels. We are able to present two examples of admissible kernels below.

The first example is Brownian bridge kernel that arises in the study of Brownian bridge stochastic process in statistics [3].

Proposition 5.1. *The Brownian bridge kernel defined by*

$$K(s, t) := \min\{s, t\} - st, \quad s, t \in (0, 1)$$

is an admissible kernel on the input space $X = (0, 1)$.

Proof. We start with validating requirement (A4). Let $0 < x_1 < x_2 < \dots < x_n < 1$ be given and $x \in (0, 1)$ be different from x_j , $j \in \mathbb{N}_n$. Direct computations show that

1. If $x < x_1$ then $K[\mathbf{x}]^{-1}K_{\mathbf{x}}(x) = \left(\frac{x}{x_1}, 0, \dots, 0\right)^T$.
2. If $x > x_n$ then $K[\mathbf{x}]^{-1}K_{\mathbf{x}}(x) = \left(0, \dots, 0, \frac{1-x}{1-x_n}\right)^T$.
3. If $x_j < x < x_{j+1}$ for some $j \in \mathbb{N}_{n-1}$ then

$$K[\mathbf{x}]^{-1}K_{\mathbf{x}}(x) = \left(0, \dots, 0, \frac{x_{j+1} - x}{x_{j+1} - x_j}, \frac{x - x_j}{x_{j+1} - x_j}, 0, \dots, 0\right)^T.$$

In all cases, it is straightforward to see $\|K[\mathbf{x}]^{-1}K_{\mathbf{x}}(x)\|_{\ell^1(\mathbb{N}_n)} \leq 1$. Therefore, requirement (A4) is indeed fulfilled.

To verify the other three requirements, we first observe

$$K(s, t) = \int_0^1 \Gamma_s(z) \Gamma_t(z) dz, \quad s, t \in (0, 1),$$

where $\Gamma_x := \chi_{(0, x)} - x$ with χ_A standing for the characteristic function of $A \subseteq (0, 1)$. Suppose that $K[\mathbf{x}]\mathbf{c} = 0$ for some $\mathbf{c} \in \mathbb{C}^n$. Then we have

$$\int_0^1 \left| \sum_{j=1}^n c_j \Gamma_{x_j}(z) \right|^2 dz = \mathbf{c}^* K[\mathbf{x}] \mathbf{c} = 0,$$

which implies that

$$\sum_{j=1}^n c_j \Gamma_{x_j}(z) = 0 \text{ for almost every } z \in [0, 1].$$

Clearly, Γ_{x_j} , $j \in \mathbb{N}_n$ are linearly independent. Therefore, $c_j = 0$ for all $j \in \mathbb{N}_n$. Requirement (A1) is hence satisfied.

The function K is clearly bounded by 1. Suppose that for some $\mathbf{c} \in \ell^1(\mathbb{N})$ and pairwise distinct $x_j \in (0, 1)$, $j \in \mathbb{N}$

$$\sum_{j=1}^{\infty} c_j K(x_j, x) = \int_0^1 \left(\sum_{j=1}^{\infty} c_j \Gamma_{x_j}(z) \right) \Gamma_x(z) dz = 0 \text{ for all } x \in (0, 1).$$

It implies that the function $\phi := \sum_{j=1}^{\infty} c_j \Gamma_{x_j}$ is orthogonal to Γ_x for all $x \in (0, 1)$, that is,

$$\int_0^x \phi(t) dt - x \int_0^1 \phi(t) dt = 0 \text{ for all } x \in (0, 1).$$

Taking the derivative on both sides of the above equations yields that ϕ equals a constant C almost everywhere on $[0, 1]$. Namely,

$$\sum_{j=1}^{\infty} c_j \chi_{[0, x_j]} - \sum_{j=1}^{\infty} c_j x_j = C \text{ almost everywhere.}$$

We now take the derivative of both sides of the equation above in the distributional sense to get $\sum_{j \in \mathbb{N}} c_j \delta_{x_j} = 0$. Let j be an arbitrary but fixed positive integer. We can find a sequence of infinitely continuously differentiable functions ϕ_k , $k \in \mathbb{N}$ such that $\|\phi_k\|_{L^\infty([0, 1])} \leq 1$, $\phi_k(x_j) = 1$, and the Lebesgue measure of the set where ϕ_k is nonzero is less than or equal to $\frac{1}{k}$. For each $N \in \mathbb{N}$, we have for sufficiently large k that

$$\phi_k(t_l) = 0 \text{ for all } l \in \mathbb{N}_N \setminus \{j\}.$$

We get for this ϕ_k

$$0 = \left| \left(\sum_{l=1}^{\infty} c_l \delta_{x_l} \right) (\phi_k) \right| \geq |c_j| - \sum_{l > N} |c_l|.$$

Since $\sum_{l > N} |c_l|$ converges to zero as $N \rightarrow \infty$, we have $c_j = 0$. Therefore, $c = 0$ for j is arbitrary chosen.

We conclude that all the four requirements of an admissible kernel are fulfilled by the Brownian bridge kernel. \square

The second example is the exponential kernel (also called the C^0 Matérn kernel).

Proposition 5.2. *The exponential kernel*

$$K(s, t) := e^{-|s-t|}, \quad s, t \in \mathbb{R} \quad (5.1)$$

is an admissible kernel on \mathbb{R} .

Proof. We have seen in Section 3 that this kernel satisfies requirements (A1)–(A3). It remains to check requirement (A4). Let $x_1 < x_2 < \dots < x_n$ be given and $x \in \mathbb{R}$ be different from x_j , $j \in \mathbb{N}_n$. Direct computations show that

1. If $x < x_1$ then $K[\mathbf{x}]^{-1} K_{\mathbf{x}}(x) = (e^{x-x_1}, 0, \dots, 0)^T$.
2. If $x > x_n$ then $K[\mathbf{x}]^{-1} K_{\mathbf{x}}(x) = (0, \dots, 0, e^{x_n-x})^T$.
3. If $x_j < x < x_{j+1}$ for some $j \in \mathbb{N}_{n-1}$ then

$$K[\mathbf{x}]^{-1} K_{\mathbf{x}}(x) = \left(0, \dots, 0, \frac{e^{x_{j+1}-x} - e^{x-x_{j+1}}}{e^{x_{j+1}-x_j} - e^{x_j-x_{j+1}}}, \frac{e^{x-x_j} - e^{x_j-x}}{e^{x_{j+1}-x_j} - e^{x_j-x_{j+1}}}, 0, \dots, 0 \right)^T.$$

In all cases, $\|K[\mathbf{x}]^{-1} K_{\mathbf{x}}(x)\|_{\ell^1(\mathbb{N}_n)} \leq 1$. The proof is complete. \square

Finally, we remark that by numerical experiments, the Gaussian kernel

$$K(s, t) = \exp\left(-\frac{(s-t)^2}{\sigma}\right), \quad s, t \in \mathbb{R}$$

does not satisfy (A4). Consequently, neither does the Gaussian kernel (3.4) on \mathbb{R}^d . The same situation happens to the inverse multiquadric (3.5) when $\beta = 1/2$.

6 Relaxation of the Admissible Condition (A4)

As seen above, the admissible condition (A4) is satisfied for few commonly used kernels. This section aims at weakening this requirement to accommodate more kernels. We are very grateful to the anonymous referee for a useful remark that inspired the approach below.

Let K be a function on $X \times X$ that satisfies (A1)-(A3) and let \mathcal{B} be constructed by (1.3). The condition (A4) is meant to ensure the validity of the linear representer theorem for regularized learning in \mathcal{B} . To see how it can be relaxed, we first examine the role of the linear representer theorem in the learning rate estimate. Consider the ℓ^1 norm coefficient-based regularization algorithm

$$\min_{\mathbf{c} \in \mathbb{C}^n} \frac{1}{n} \sum_{j=1}^n |K^{\mathbf{x}}(x_j) \mathbf{c} - y_j|^2 + \mu \|\mathbf{c}\|_{\ell^1(\mathbb{N}_n)} \quad (6.1)$$

where $\mathbf{x} := \{x_j : j \in \mathbb{N}_n\}$ is a sequence of sampling points from the input space X , $y_j \in Y \subseteq \mathbb{C}$ is the observed output on x_j , μ is a positive regularization parameter. Following a commonly used assumption in machine learning, we assume that the sample data $\mathbf{z} := \{(x_j, y_j) : j \in \mathbb{N}_n\} \in X \times Y$ is formed by independent and identically distributed instances of a random variable $(x, y) \in X \times Y$ subject to an unknown probability measure ρ on $X \times Y$. Let $\mathbf{c}_{\mathbf{z}, \mu}$ be a minimizer of (6.1). We hope that the obtained function

$$f_{\mathbf{z}, \mu}(x) := K^{\mathbf{x}}(x) \mathbf{c}_{\mathbf{z}, \mu}, \quad x \in X \quad (6.2)$$

will well predict the outputs of new inputs from X . The performance of a general predictor $f : X \rightarrow Y$ is usually measured by

$$\mathcal{E}(f) := \int_{X \times Y} |f(x) - y|^2 d\rho.$$

The predictor that minimizes the above error is the regression function

$$f_\rho(x) := \int_Y y d\rho(y|x), \quad x \in X,$$

where $\rho(y|x)$ denotes the conditional probability measure of y with respect to x . This optimal predictor f_ρ is unreachable as ρ is unknown. We shall approximate f_ρ with $f_{\mathbf{z}, \mu}$. More precisely, we expect with a large confidence that the approximation error $\mathcal{E}(f_{\mathbf{z}, \mu}) - \mathcal{E}(f_\rho)$ would converge to zero fast as the number of sampling points increases.

A standard approach [7] in estimating the error $\mathcal{E}(f_{\mathbf{z}, \mu}) - \mathcal{E}(f_\rho)$ is to bound it by the sum of the sampling error, the hypothesis error and the regularization error. Let g be an arbitrary function from \mathcal{B} and set for each function $f : X \rightarrow \mathbb{C}$

$$\mathcal{E}_{\mathbf{z}}(f) := \frac{1}{n} \sum_{j=1}^n |f(x_j) - y_j|^2.$$

The approximation error $\mathcal{E}(f_{\mathbf{z}, \mu}) - \mathcal{E}(f_\rho)$ can then be decomposed into the sum of four quantities

$$\mathcal{E}(f_{\mathbf{z}, \mu}) - \mathcal{E}(f_\rho) = \mathcal{S}(\mathbf{z}, \mu, g) + \mathcal{P}(\mathbf{z}, \mu, g) + \mathcal{D}(\mu, g) - \mu \|f_{\mathbf{z}, \mu}\|_{\mathcal{B}},$$

where the *sampling error*, the *hypothesis error* and the *regularization error* are respectively defined by

$$\begin{aligned} \mathcal{S}(\mathbf{z}, \mu, g) &:= \mathcal{E}(f_{\mathbf{z}, \mu}) - \mathcal{E}_{\mathbf{z}}(f_{\mathbf{z}, \mu}) + \mathcal{E}_{\mathbf{z}}(g) - \mathcal{E}(g), \\ \mathcal{P}(\mathbf{z}, \mu, g) &:= (\mathcal{E}_{\mathbf{z}}(f_{\mathbf{z}, \mu}) + \mu \|f_{\mathbf{z}, \mu}\|_{\mathcal{B}}) - (\mathcal{E}_{\mathbf{z}}(g) + \mu \|g\|_{\mathcal{B}}), \\ \mathcal{D}(\mu, g) &:= \mathcal{E}(g) - \mathcal{E}(f_\rho) + \mu \|g\|_{\mathcal{B}}. \end{aligned}$$

Under the condition (A4), \mathcal{B} satisfies the linear representer theorem. As a result,

$$\mathcal{E}_{\mathbf{z}}(f_{\mathbf{z},\mu}) + \mu\|f_{\mathbf{z},\mu}\|_{\mathcal{B}} = \min_{f \in \mathcal{S}^{\mathbf{x}}} \mathcal{E}_{\mathbf{z}}(f) + \mu\|f\|_{\mathcal{B}} = \min_{f \in \mathcal{B}} \mathcal{E}_{\mathbf{z}}(f) + \mu\|f\|_{\mathcal{B}}. \quad (6.3)$$

Immediately, one has $\mathcal{P}(\mathbf{z}, \mu, g) \leq 0$, leading to the estimate

$$\mathcal{E}(f_{\mathbf{z},\mu}) - \mathcal{E}(f_{\rho}) \leq \mathcal{S}(\mathbf{z}, \mu, g) + \mathcal{D}(\mu, g).$$

Starting from the above inequality, learning rates of $f_{\mathbf{z},\mu}$ can be obtained [25]. To weaken (A4), we should not stick to the linear representer theorem (6.3). Instead, we wish to replace it with the *relaxed linear representer theorem*

$$\min_{f \in \mathcal{S}^{\mathbf{x}}} \mathcal{E}_{\mathbf{z}}(f) + \mu\|f\|_{\mathcal{B}} \leq \min_{f \in \mathcal{B}} \mathcal{E}_{\mathbf{z}}(f) + \mu\beta_n\|f\|_{\mathcal{B}}, \quad (6.4)$$

where β_n is a constant depending on the number n of sampling points, the kernel K and the input space X . For simplicity, we suppress the notations K and X as they are fixed in our context. The approximation error $\mathcal{E}(f_{\mathbf{z},\mu}) - \mathcal{E}(f_{\rho})$ is accordingly factored as

$$\mathcal{E}(f_{\mathbf{z},\mu}) - \mathcal{E}(f_{\rho}) = \mathcal{S}(\mathbf{z}, \mu, g) + \tilde{\mathcal{P}}(\mathbf{z}, \mu, g) + \tilde{\mathcal{D}}(\mu, g) - \mu\|f_{\mathbf{z},\mu}\|_{\mathcal{B}},$$

where

$$\begin{aligned} \tilde{\mathcal{P}}(\mathbf{z}, \mu, g) &:= (\mathcal{E}_{\mathbf{z}}(f_{\mathbf{z},\mu}) + \mu\|f_{\mathbf{z},\mu}\|_{\mathcal{B}}) - (\mathcal{E}_{\mathbf{z}}(g) + \mu\beta_n\|g\|_{\mathcal{B}}), \\ \tilde{\mathcal{D}}(\mu, g) &:= \mathcal{E}(g) - \mathcal{E}(f_{\rho}) + \mu\beta_n\|g\|_{\mathcal{B}}. \end{aligned}$$

By (6.4), we keep the advantage that $\tilde{\mathcal{P}}(\mathbf{z}, \mu, g) \leq 0$. Therefore,

$$\mathcal{E}(f_{\mathbf{z},\mu}) - \mathcal{E}(f_{\rho}) \leq \mathcal{S}(\mathbf{z}, \mu, g) + \tilde{\mathcal{D}}(\mu, g).$$

As long as β_n does not increase too fast as n increases, one is still able to obtain a learning rate competitive with those in [25, 30]. We shall omit the detailed arguments and assumptions on the kernel K , the regression function f_{ρ} and the input space X , as they are similar to those in [25]. We present one result that for all $0 < \delta < 1$, there exists a constant C_{δ} such that with confidence $1 - \delta$, we have

$$\mathcal{E}(f_{\mathbf{z},\mu}) - \mathcal{E}(f_{\rho}) \leq C_{\delta} \left((\mu\beta_n)^{\frac{2s}{1+s}} + \frac{\log \frac{2}{\delta}}{n} (\mu\beta_n)^{\frac{2s-2}{1+s}} + \frac{\log \frac{2}{\delta}}{\sqrt{n}} (\mu\beta_n)^{\frac{2s-1}{1+s}} + \frac{\log \frac{2}{\delta} + \log(1+n)}{(\mu\beta_n)^2} \beta_n^2 n^{-\frac{1}{1+\theta}} \right),$$

where $s \in (0, 1)$ represents the regularity of f_{ρ} , $\theta > 0$ is a positive constant related to assumptions on the kernel K and the input space X , [25]. Thus, as long as β_n^2 does not cancel the decay of the term $n^{-\frac{1}{1+\theta}}$, one still has the hope of getting a satisfactory learning rate when μ is appropriately chosen. We discuss two instances below:

(i) If β_n is uniformly bounded with a large confidence then $\mathcal{E}(f_{\mathbf{z},\mu}) - \mathcal{E}(f_{\rho})$ has the same learning rate as that established in [25], that is,

$$\mathcal{E}(f_{\mathbf{z},\mu}) - \mathcal{E}(f_{\rho}) \leq C_{\delta} n^{-\frac{s}{1+2s} \frac{1}{1+\theta}} \log \frac{2+2n}{\delta}. \quad (6.5)$$

(ii) If $\beta_n \leq Cn^{\alpha}$ for some positive constants C and $\alpha < \frac{1}{2+2\theta}$ then

$$\mathcal{E}(f_{\mathbf{z},\mu}) - \mathcal{E}(f_{\rho}) \leq C_{\delta} n^{-\frac{s}{1+2s} (\frac{1}{1+\theta} - 2\alpha)} \log \frac{2+2n}{\delta}. \quad (6.6)$$

If we give up the linear representer theorem and pursue the relaxed version (6.4) instead, how can the admissible condition (A4) be weakened? We next answer this question.

Proposition 6.1. *If there exists some $\beta_n \geq 1$ such that for all $\mathbf{y} \in \mathbb{C}^n$*

$$\min_{f \in \mathcal{I}_{\mathbf{x}}(\mathbf{y})} \|f\|_{\mathcal{B}} \geq \frac{1}{\beta_n} \min_{\mathcal{I}_{\mathbf{x}}(\mathbf{y}) \cap \mathcal{S}^{\mathbf{x}}} \|f\|_{\mathcal{B}} \quad (6.7)$$

then the relaxed linear representer theorem (6.4) holds true for any continuous loss function V and any regularization parameter μ .

Proof. Suppose that (6.7) is satisfied. Let f_0 be a minimizer of

$$\min_{f \in \mathcal{B}} V(f(\mathbf{x})) + \lambda \beta_n \|f\|_{\mathcal{B}}.$$

Choose g to be a function in $\mathcal{S}^{\mathbf{x}}$ that interpolates f_0 at \mathbf{x} , namely, $g(\mathbf{x}) = f_0(\mathbf{x})$. By (6.7),

$$\|g\|_{\mathcal{B}} \leq \beta_n \|f_0\|_{\mathcal{B}},$$

which yields

$$V(g(\mathbf{x})) + \lambda \|g\|_{\mathcal{B}} \leq V(f_0(\mathbf{x})) + \lambda \beta_n \|f_0\|_{\mathcal{B}}.$$

The proof is hence complete. \square

We next give a characterization of (6.7), which gives rise to a relaxation of the admissible condition (A4) and leads to the relaxed linear representer theorem (6.4).

Theorem 6.2. *Equation (6.7) holds true for all $\mathbf{y} \in \mathbb{C}^n$ if and only if*

$$\|(K[\mathbf{x}])^{-1} K_{\mathbf{x}}(t)\|_{\ell^1(\mathbb{N}_n)} \leq \beta_n \text{ for all } t \in X. \quad (6.8)$$

Proof. The set $\mathcal{I}_{\mathbf{x}}(\mathbf{y}) \cap \mathcal{S}^{\mathbf{x}}$ consists of only one function $f_0 := K^{\mathbf{x}}(\cdot) K[\mathbf{x}]^{-1} \mathbf{y}$. Let g be an arbitrary function in $\mathcal{I}_{\mathbf{x}}(\mathbf{y}) \cap \mathcal{B}_0$. By adding sampling points and assigning the corresponding coefficients to be zero if necessary, we may assume $g \in \mathcal{S}^{\mathbf{x} \cup \mathbf{t}} \cap \mathcal{I}_{\mathbf{x}}(\mathbf{y})$ for some $\mathbf{t} := \{t_j \in X : j \in \mathbb{N}_m\}$ disjoint with \mathbf{x} . Let $\mathbf{b} := g(\mathbf{t})$, and denote by $K[\mathbf{t}, \mathbf{x}]$ and $K[\mathbf{x}, \mathbf{t}]$ the $n \times m$ and $m \times n$ matrices given by

$$(K[\mathbf{t}, \mathbf{x}])_{jk} := K(t_k, x_j), \quad j \in \mathbb{N}_n, k \in \mathbb{N}_m, \quad (K[\mathbf{x}, \mathbf{t}])_{jk} := K(x_k, t_j) : \quad j \in \mathbb{N}_m, k \in \mathbb{N}_n.$$

Then

$$\|g\|_{\mathcal{B}} = \left\| \begin{pmatrix} K[\mathbf{x}] & K[\mathbf{t}, \mathbf{x}] \\ K[\mathbf{x}, \mathbf{t}] & K[\mathbf{t}] \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{y} \\ \mathbf{b} \end{pmatrix} \right\|_{\ell^1(\mathbb{N}_{n+m})} = \left\| \begin{pmatrix} K[\mathbf{x}]^{-1} \mathbf{y} - K[\mathbf{x}]^{-1} K[\mathbf{t}, \mathbf{x}] \tilde{\mathbf{b}} \\ \tilde{\mathbf{b}} \end{pmatrix} \right\|_{\ell^1(\mathbb{N}_{n+m})}, \quad (6.9)$$

where

$$\tilde{\mathbf{b}} := (K[\mathbf{t}] - K[\mathbf{x}, \mathbf{t}] K[\mathbf{x}]^{-1} K[\mathbf{t}, \mathbf{x}])^{-1} (\mathbf{b} - K[\mathbf{x}, \mathbf{t}] K[\mathbf{x}]^{-1} \mathbf{y}).$$

Note that as \mathbf{b} is allowed to equal any vector in \mathbb{C}^m , so is $\tilde{\mathbf{b}}$.

If (6.7) holds true for all $\mathbf{y} \in \mathbb{C}^n$ then we choose \mathbf{t} to be a singleton $\{t\}$, $\tilde{b} = 1$, and $\mathbf{y} = K[t, \mathbf{x}] = K_{\mathbf{x}}(t)$ to get

$$\left\| \begin{pmatrix} \mathbf{0} \\ 1 \end{pmatrix} \right\|_{\ell^1(\mathbb{N}_{n+1})} \geq \frac{1}{\beta_n} \|f_0\|_{\mathcal{B}} = \frac{1}{\beta_n} \|K[\mathbf{x}]^{-1} \mathbf{y}\|_{\ell^1(\mathbb{N}_n)} = \frac{1}{\beta_n} \|K[\mathbf{x}]^{-1} K_{\mathbf{x}}(t)\|_{\ell^1(\mathbb{N}_n)},$$

which is (6.8). Conversely, suppose that (6.8) is satisfied. We need to show that for all $g \in \mathcal{I}_{\mathbf{x}}(\mathbf{y})$

$$\|g\|_{\mathcal{B}} \geq \frac{1}{\beta_n} \|f_0\|_{\mathcal{B}} = \frac{1}{\beta_n} \|K[\mathbf{x}]^{-1} \mathbf{y}\|_{\ell^1(\mathbb{N}_n)}.$$

We shall discuss the case when $g \in \mathcal{I}_{\mathbf{x}}(\mathbf{y}) \cap \mathcal{B}_0$ only as the general case will then follow by the same arguments as those in the last paragraph of the proof of Theorem 4.8. Let $g \in \mathcal{I}_{\mathbf{x}}(\mathbf{y}) \cap \mathcal{B}_0$ have the norm (6.9). Clearly,

$$\|g\|_{\mathcal{B}} \geq \frac{1}{\beta_n} \|K[\mathbf{x}]^{-1} \mathbf{y}\|_{\ell^1(\mathbb{N}_n)}$$

if $\|K[\mathbf{x}]^{-1} \mathbf{y}\|_{\ell^1(\mathbb{N}_n)} \leq \beta_n \|\tilde{\mathbf{b}}\|_{\ell^1(\mathbb{N}_m)}$. When $\|K[\mathbf{x}]^{-1} \mathbf{y}\|_{\ell^1(\mathbb{N}_n)} > \beta_n \|\tilde{\mathbf{b}}\|_{\ell^1(\mathbb{N}_m)}$, we have

$$\begin{aligned} \|g\|_{\mathcal{B}} &\geq \|K[\mathbf{x}]^{-1} \mathbf{y}\|_{\ell^1(\mathbb{N}_n)} - \|K[\mathbf{x}]^{-1} K[\mathbf{t}, \mathbf{x}] \tilde{\mathbf{b}}\|_{\ell^1(\mathbb{N}_n)} + \|\tilde{\mathbf{b}}\|_{\ell^1(\mathbb{N}_m)} \\ &\geq \|K[\mathbf{x}]^{-1} \mathbf{y}\|_{\ell^1(\mathbb{N}_n)} - \left(\max_{k \in \mathbb{N}_m} \|K[\mathbf{x}]^{-1} K_{\mathbf{x}}(t_k)\|_{\ell^1(\mathbb{N}_n)} \right) \|\tilde{\mathbf{b}}\|_{\ell^1(\mathbb{N}_m)} + \|\tilde{\mathbf{b}}\|_{\ell^1(\mathbb{N}_m)} \\ &\geq \|K[\mathbf{x}]^{-1} \mathbf{y}\|_{\ell^1(\mathbb{N}_n)} - (\beta_n - 1) \|\tilde{\mathbf{b}}\|_{\ell^1(\mathbb{N}_m)} \geq \|K[\mathbf{x}]^{-1} \mathbf{y}\|_{\ell^1(\mathbb{N}_n)} - (\beta_n - 1) \frac{1}{\beta_n} \|K[\mathbf{x}]^{-1} \mathbf{y}\|_{\ell^1(\mathbb{N}_n)} \\ &= \frac{1}{\beta_n} \|K[\mathbf{x}]^{-1} \mathbf{y}\|_{\ell^1(\mathbb{N}_n)}, \end{aligned}$$

which completes the proof. \square

The above result together with the discussion of the application of Proposition 6.1 to regularized learning provides a relaxation of the requirement (A4). The quantity $\sup_{t \in X} \|K[\mathbf{x}]^{-1} K_{\mathbf{x}}(t)\|_{\ell^1(\mathbb{N}_n)}$ is the Lebesgue constant of the kernel interpolation. Asking it to be exactly bounded by 1 is indeed demanding. Recent numerical experiments [8] and analysis [12] indicate that for many kernels, this Lebesgue constant could be uniformly bounded. In this case, the ℓ^1 -regularized learning in \mathcal{B} performs well by (6.5). Furthermore, as long as β_n does not increase to infinity too fast, the learning scheme can still work well by (6.6). Specifically, it was proved in [12] that the Lebesgue constant for the reproducing kernel of the Sobolev space on a compact domain is uniformly bounded for quasi-uniform input points (see, Theorem 4.6 therein). Another example is given in [8] for translation invariant kernels $K(x, y) = \phi(x - y)$, $x, y \in \mathbb{R}^d$. It was shown there that as long as

$$c_1(1 + \|\xi\|_2^2)^{-\tau} \leq \hat{\phi}(\xi) \leq c_2(1 + \|\xi\|_2^2)^{-\tau}, \quad \|\xi\|_2 > M \quad (6.10)$$

for some positive constants c_1, c_2, M and τ , the Lebesgue constant for quasi-uniform inputs is bounded by a multiple of \sqrt{n} . Commonly used kernels satisfying (6.10) include Poisson radial functions [10], Matérn kernels and Wendland's compactly supported kernels [28]. Finally, we remark from numerical experiments that the following kernels [20]

$$\exp\left(-\|x - y\|_{\ell^p(\mathbb{N}_d)}^\gamma\right), \quad x, y \in \mathbb{R}^d, \quad \gamma \in (0, 1), \quad p = 1, 2$$

seem to satisfy (A4) for small enough γ and moderate n . We shall leave the search of more kernels satisfying (A4) and its relaxation (6.8) as an open question for future study.

7 Numerical Experiments

We end this paper with a numerical experiment to show that the regularization algorithm (4.1) is indeed able to yield sparse learning compared to the classical regularization network in machine learning.

We shall use the exponential kernel K (5.1). Let \mathcal{B} be the corresponding RKBS with the ℓ^1 norm constructed by (1.3) and let \mathcal{H}_K be the RKHS of K . We restrict ourselves to the field of real numbers and use the square loss function $V(f(\mathbf{x})) := \|f(\mathbf{x}) - \mathbf{y}\|_2^2$. We shall compare the two models

$$\min_{f \in \mathcal{B}} \|f(\mathbf{x}) - \mathbf{y}\|_2^2 + \mu \|f\|_{\mathcal{B}}$$

and

$$\min_{g \in \mathcal{H}_K} \|g(\mathbf{x}) - \mathbf{y}\|_2^2 + \mu \|g\|_{\mathcal{H}_K}^2.$$

Both of them satisfy the linear representer theorem. Specifically, the minimizers f_0 and g_0 of the above two models are respectively given by

$$f_0 = K^{\mathbf{x}}(\cdot) \mathbf{b} \text{ with } \mathbf{b} := \operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^n} \{\|K[\mathbf{x}]\mathbf{c} - \mathbf{y}\|_2^2 + \mu \|\mathbf{c}\|_{\ell^1(\mathbb{N}_n)}\}$$

and

$$g_0 = K^{\mathbf{x}}(\cdot) \mathbf{h} \text{ with } \mathbf{h} := \operatorname{argmin}_{\mathbf{c} \in \mathbb{R}^n} \{\|K[\mathbf{x}]\mathbf{c} - \mathbf{y}\|_2^2 + \mu \mathbf{c}^T K[\mathbf{x}]\mathbf{c}\}.$$

We point out that the above ℓ^1 minimization problem about \mathbf{b} does not have a closed form solution. There are numerous methods proposed to solve this problem and here we employ the proximity algorithm recently developed in [18]. The closed form of the minimizer \mathbf{h} is well known to be $(K[\mathbf{x}] + \mu I_n)^{-1} \mathbf{y}$. Here I_n denotes the $n \times n$ identity matrix.

For both models, \mathbf{x} is set to be 200 equally spaced points in $[-1, 1]$ and the output vector \mathbf{y} is chosen to be the evaluation of the target function

$$f(x) = e^{-|x+1|} + e^{-|x+0.8|} + e^{-|x|} + e^{-|x-0.8|} + e^{-|x-1|}, \quad x \in [-1, 1]$$

at \mathbf{x} and then disturbed by some noise. Also, the regularization parameter μ for each model will be optimally chosen from $\{10^j : j = -7, -6, \dots, 1\}$ so that the distance between the learned function and the target function in $L^2([-1, 1])$ will be minimized. We then compare the approximation accuracy measured by this error and the sparsity for these two models. The sparsity is measured by the number of nonzero components in the coefficient vectors \mathbf{b} and \mathbf{h} .

	Gaussian noise		Uniform noise		Pepper sauce noise	
	Error	Sparsity (Max)	Error	Sparsity (Max)	Error	Sparsity (Max)
RKHS	2.1E-3	200 (200)	7.9E-4	200 (200)	9.4E-4	200 (200)
RKBS	1.0E-3	13.4 (17)	3.6E-4	14.7 (25)	4.5E-4	14.5 (23)

Table 1: Comparison of the least square regularization in RKHS and in RKBS with the ℓ^1 norm for the exponential kernel.

We test both models with three types of noise: Gaussian noise with variance 0.01, uniform noise in $[-0.1, 0.1]$ and some random pepper sauce noise in $\{-0.1, 0.1\}$. For each type of noise, we run 50 times of numerical experiments and compute the average approximation error, the average sparsity, and the maximum sparsity in the 50 experiments. The results are tabulated above.

References

- [1] A. Argyriou, C. A. Micchelli, and M. Pontil. When is there a representer theorem? Vector versus matrix regularizers. *J. Mach. Learn. Res.*, 10:2507–2529, 2009.
- [2] N. Aronszajn. Theory of reproducing kernels. *Trans. Amer. Math. Soc.*, 68:337–404, 1950.
- [3] A. Berlinet and C. Thomas-Agnan. *Reproducing Kernel Hilbert Spaces in Probability and Statistics*. Kluwer, Dordrecht, 2004.
- [4] E. J. Candès, J. Romberg, and T. Tao. Robust uncertainty principles: exact signal reconstruction from highly incomplete frequency information. *IEEE Trans. Inform. Theory*, 52(2):489–509, 2006.
- [5] S. S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. *SIAM J. Sci. Comput.*, 20(1):33–61, 1998.
- [6] F. Cucker and S. Smale. On the mathematical foundations of learning. *Bull. Amer. Math. Soc. (N.S.)*, 39(1):1–49 (electronic), 2002.
- [7] F. Cucker and D.-X. Zhou. *Learning theory: an approximation theory viewpoint*. Cambridge Monographs on Applied and Computational Mathematics. Cambridge University Press, Cambridge, 2007. With a foreword by Stephen Smale.
- [8] S. De Marchi and R. Schaback. Stability of kernel-based interpolation. *Adv. Comput. Math.*, 32(2):155–161, 2010.
- [9] T. Evgeniou, M. Pontil, and T. Poggio. Regularization networks and support vector machines. *Adv. Comput. Math.*, 13(1):1–50, 2000.
- [10] B. Fornberg, E. Larsson, and G. Wright. A new class of oscillatory radial basis functions. *Comput. Math. Appl.*, 51(8):1209–1222, 2006.
- [11] J. R. Giles. Classes of semi-inner-product spaces. *Trans. Amer. Math. Soc.*, 129:436–446, 1967.
- [12] T. Hangelsbroek, F. J. Narcowich, and J. D. Ward. Kernel approximation on manifolds I: bounding the Lebesgue constant. *SIAM J. Math. Anal.*, 42(4):1732–1760, 2010.
- [13] R. C. James. Characterizations of reflexivity. *Studia Math.*, 23:205–216, 1963/1964.
- [14] G. Kimeldorf and G. Wahba. Some results on Tchebycheffian spline functions. *J. Math. Anal. Appl.*, 33:82–95, 1971.
- [15] G. Lumer. Semi-inner-product spaces. *Trans. Amer. Math. Soc.*, 100:29–43, 1961.
- [16] C. A. Micchelli and A. Pinkus. Variational problems arising from balancing several error criteria. *Rendiconti di Matematica, Serie VII*, 14:37–86, 1994.
- [17] C. A. Micchelli and M. Pontil. On learning vector-valued functions. *Neural Comput.*, 17(1):177–204, 2005.
- [18] C. A. Micchelli, L. Shen, and Y. Xu. Proximity algorithms for image models: denoising. *Inverse Problems*, 27:045009, 2011.

- [19] C. A. Micchelli, Y. Xu, and H. Zhang. Universal kernels. *J. Mach. Learn. Res.*, 7:2651–2667, 2006.
- [20] I. J. Schoenberg. Metric spaces and positive definite functions. *Trans. Amer. Math. Soc.*, 44(3):522–536, 1938.
- [21] B. Schölkopf, R. Herbrich, and A. J. Smola. A generalized representer theorem. In *Computational learning theory (Amsterdam, 2001)*, volume 2111 of *Lecture Notes in Comput. Sci.*, pages 416–426. Springer, Berlin, 2001.
- [22] B. Schölkopf and A. J. Smola. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond (Adaptive Computation and Machine Learning)*. The MIT Press, Cambridge, December 2001.
- [23] J. Shawe-Taylor and N. Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge University Press, Cambridge, 2004.
- [24] G. Song and Y. Xu. Approximation of high-dimensional kernel matrices by multilevel circulant matrices. *J. Complexity*, 26(4):375–405, 2010.
- [25] G. Song and H. Zhang. Reproducing kernel banach spaces with the ℓ^1 norm ii: error analysis for regularized least square regression. *Neural Comput.*, 23(10):2713–2729, 2011.
- [26] R. Tibshirani. Regression shrinkage and selection via the lasso. *J. Roy. Statist. Soc. Ser. B*, 58(1):267–288, 1996.
- [27] V. N. Vapnik. *Statistical Learning Theory*. Adaptive and Learning Systems for Signal Processing, Communications, and Control. John Wiley & Sons Inc., New York, 1998. A Wiley-Interscience Publication.
- [28] H. Wendland. *Scattered data approximation*, volume 17 of *Cambridge Monographs on Applied and Computational Mathematics*. Cambridge University Press, Cambridge, 2005.
- [29] Z. M. Wu. Compactly supported positive definite radial functions. *Adv. Comput. Math.*, 4(3):283–292, 1995.
- [30] Q.-W. Xiao and D.-X. Zhou. Learning by nonsymmetric kernels with data dependent spaces and ℓ^1 -regularizer. *Taiwanese J. Math.*, 14(5):1821–1836, 2010.
- [31] H. Zhang, Y. Xu, and J. Zhang. Reproducing kernel Banach spaces for machine learning. *J. Mach. Learn. Res.*, 10:2741–2775, 2009.
- [32] H. Zhang and J. Zhang. Regularized learning in Banach spaces as an optimization problem: representer theorems. *J. Global Optim.* to appear.
- [33] H. Zhang and J. Zhang. Frames, Riesz bases, and sampling expansions in Banach spaces via semi-inner products. *Appl. Comput. Harmon. Anal.*, 31:1–25, 2011.